

*Module*

**17**

ECE M216A

# Packaging

**Prof. Dejan Marković**

ee216a@gmail.com

# Package Functions

---

- ② • **Mechanical connection of chip to board**
- ① • **Electrical connection of signals and power**
  - Short wires with low R and L
- ③ • **Removes heat produced on chip**
  - **Protects chip from mechanical damage**

[Reading: Weste, Harris VLSI book]

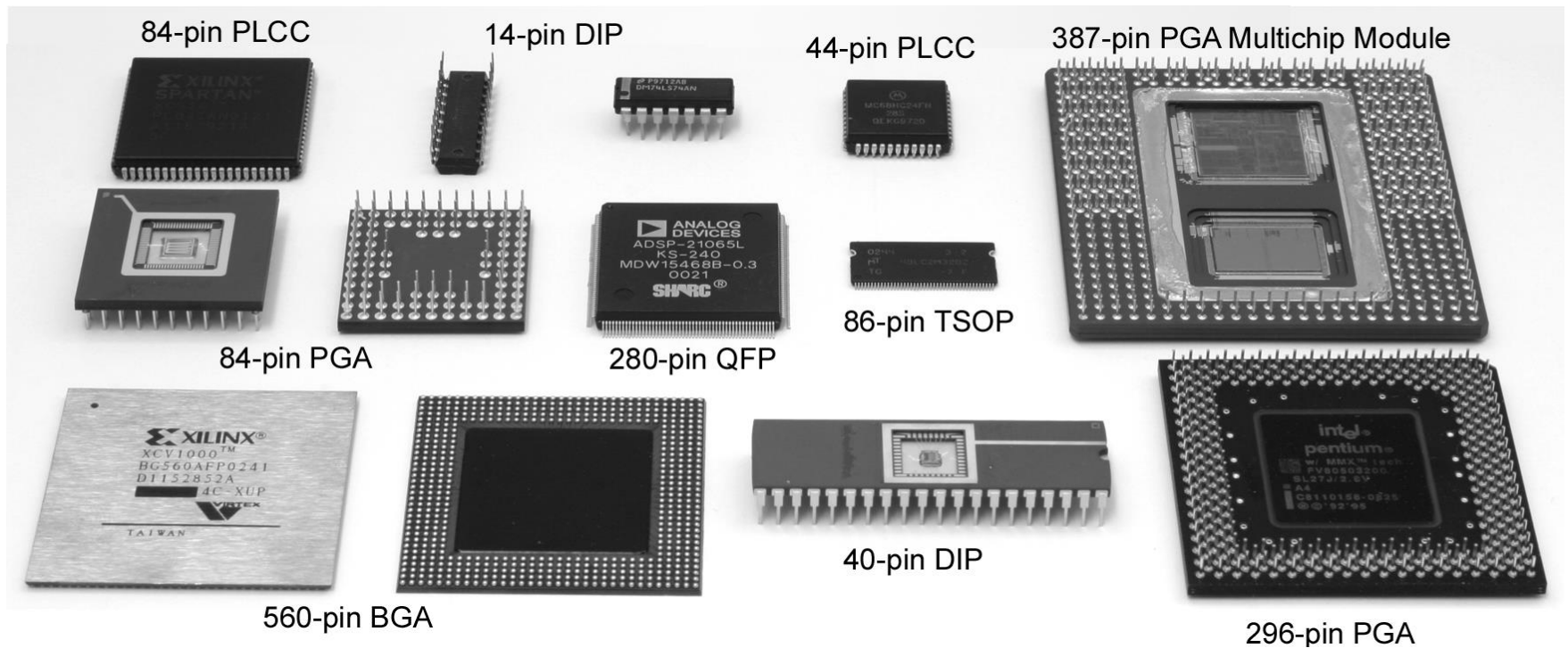
# Main Issues

---

- **Cost**
- **Thermal impedance:**  
how effectively package removes heat from the die
- **Lead inductance**
  - Ceramic pin grid array package – lowest
  - Cheap epoxy plastic – highest

# Basic Package Types

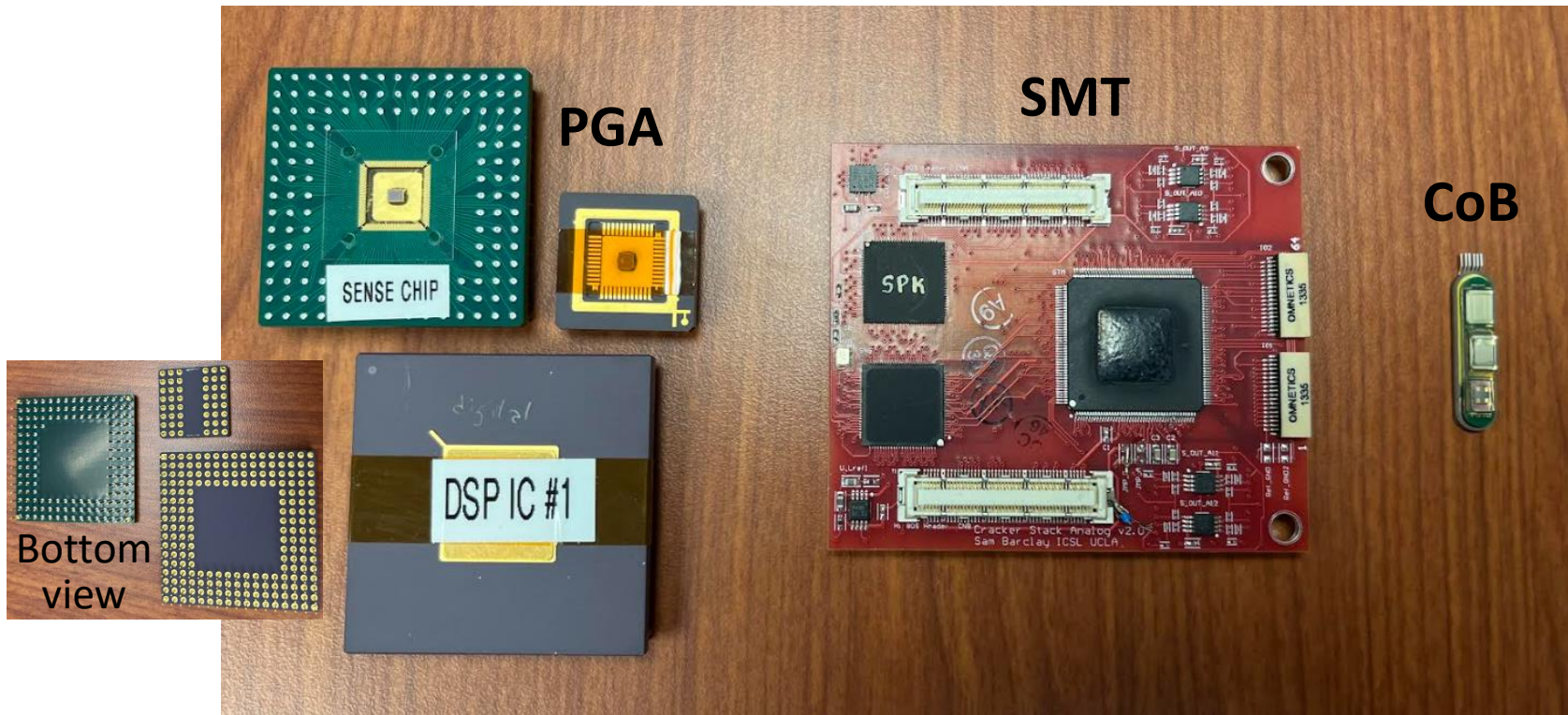
- **PLCC:** plastic leadless chip carrier
- **TSOP:** thin small outline package
- **QFP:** quad flat pack
- **DIP:** dual inline package
- **PGA:** pin grid array
- **BGA:** ball grid array



[Weste, Harris VLSI book, Fig. 12.1]

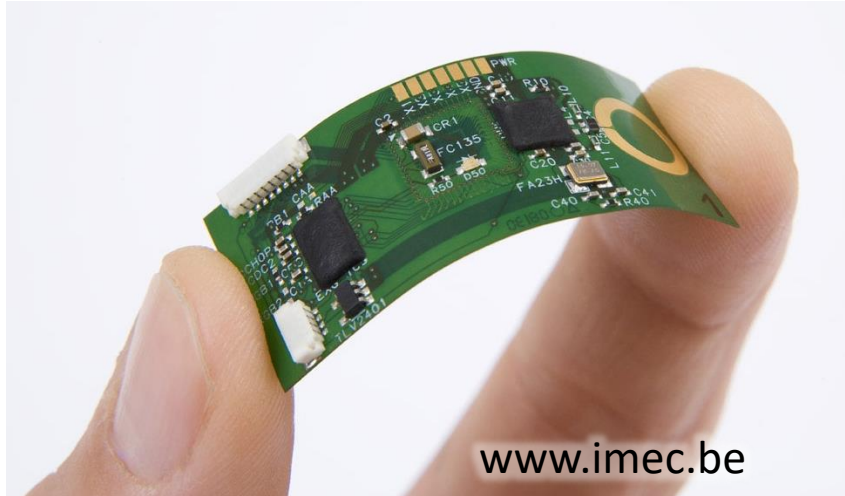
# Purpose-Driven Packaging (Same Chipset)

- Packages: PGA (left), SMT (middle) and CoB (right) for different needs / use cases



# Advanced Packaging

## Thin-chip



- Chip thinned to 25 $\mu$ m
- Flexible ultra-thin package
- Embedded in flex PCB

## System-in-package (SiP)

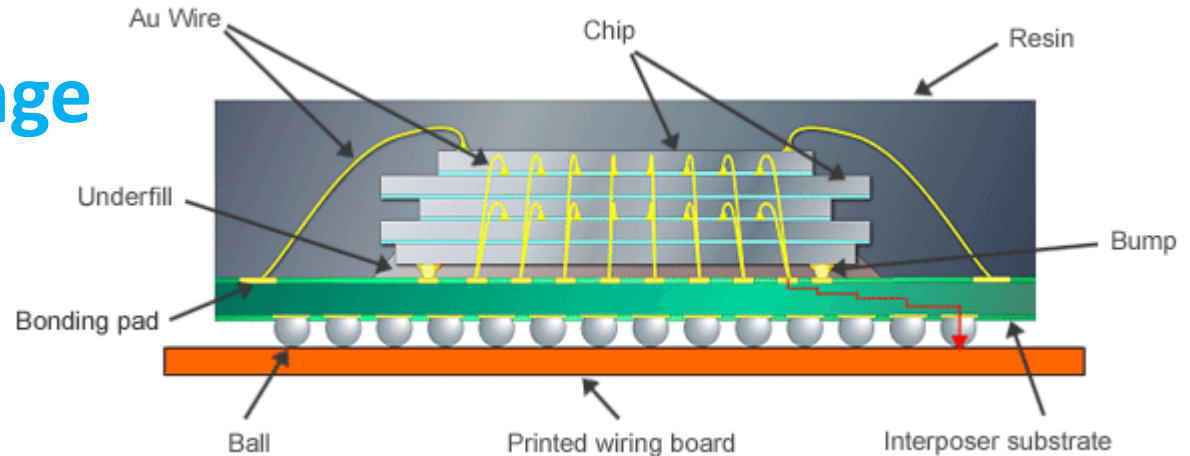
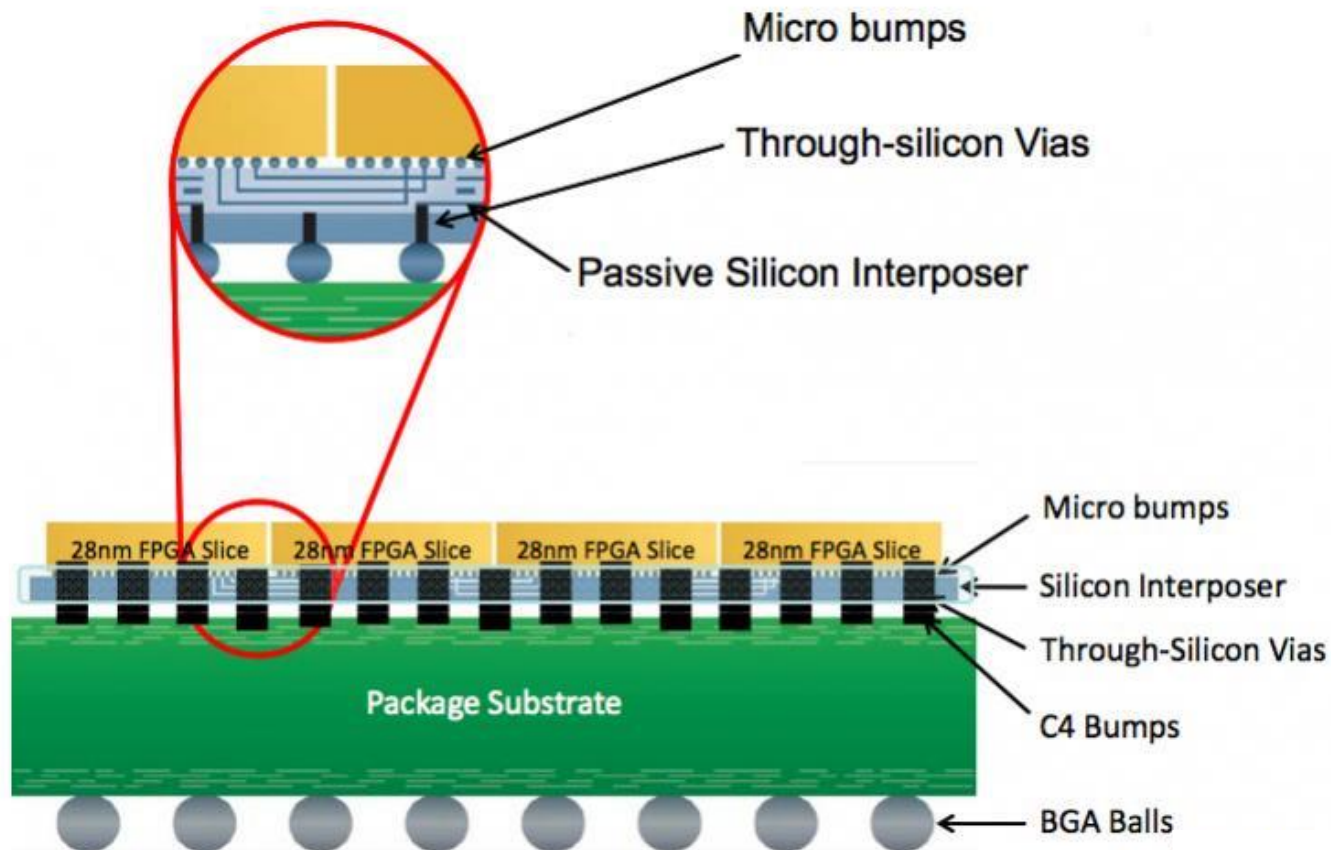


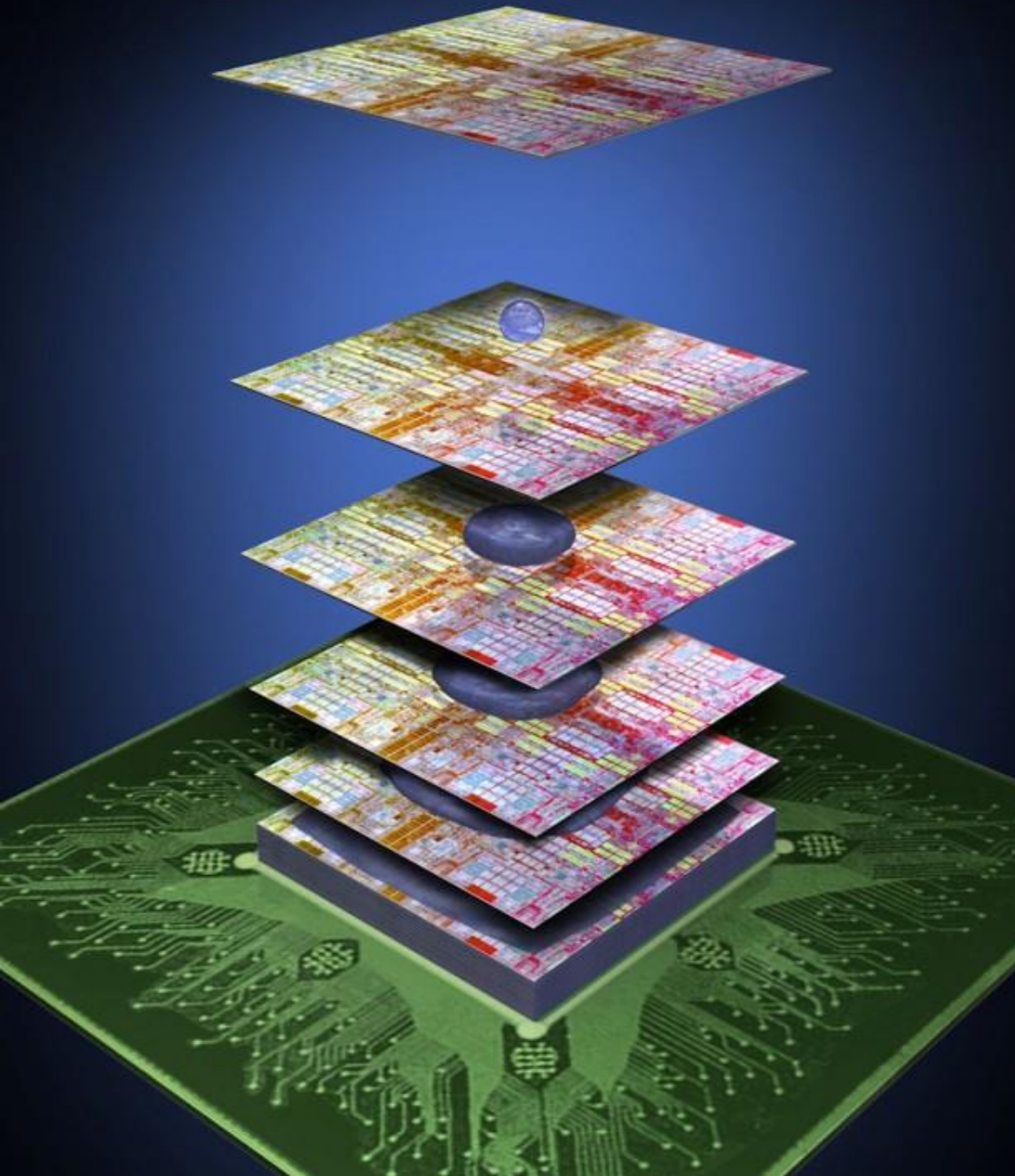
Photo: [www.renesas.com](http://www.renesas.com)

# Example: Xilinx Virtex 7 Package

- Controlled collapse chip connection (C4) bumps
- Through-silicon vias (TSVs)



# Si Skyscraper



- **3M** and **IBM**  
(2011)
- 3M glue
- Up to 100 chips
- Goal: bond wafers

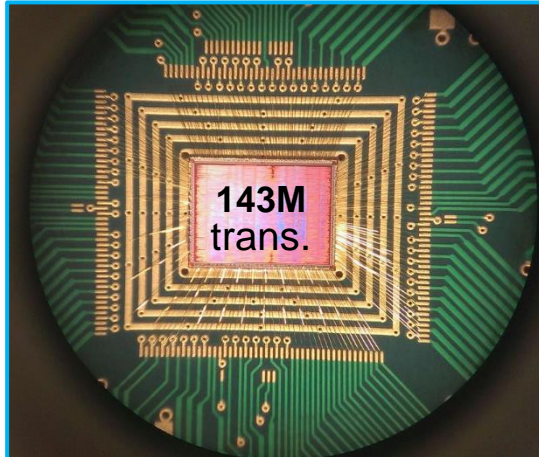
<http://youtu.be/rbj5vrXulD0>

Photo: IBM



# Example Chip(let) Assemblies

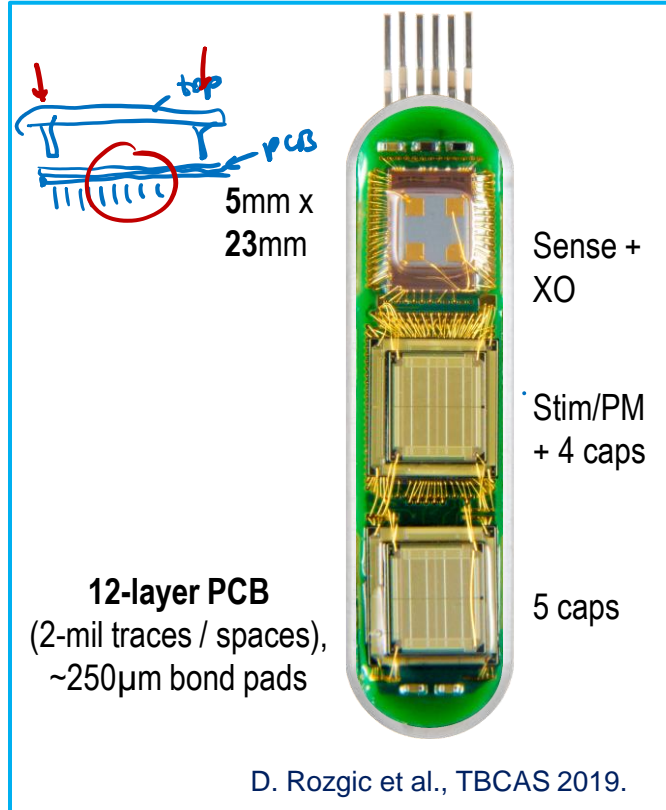
Courtesy of DMGroup



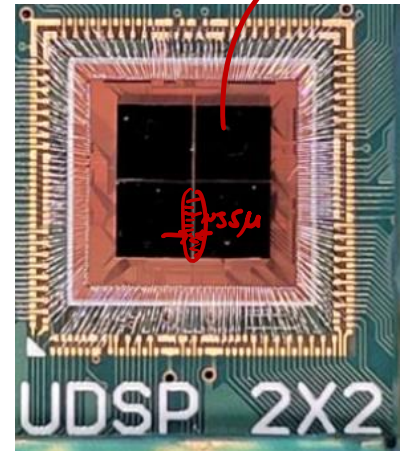
C. Wang et al., ISSCC 2014.

Slice L/M
Slice L/M
Slice L/M
Slice L
DSP-48, Slice L, BRAM
Slice L/M
Slice L/M
64-8k FFT
16-core UDSP
Slice L/M

24.5mm<sup>2</sup>  
(40nm)



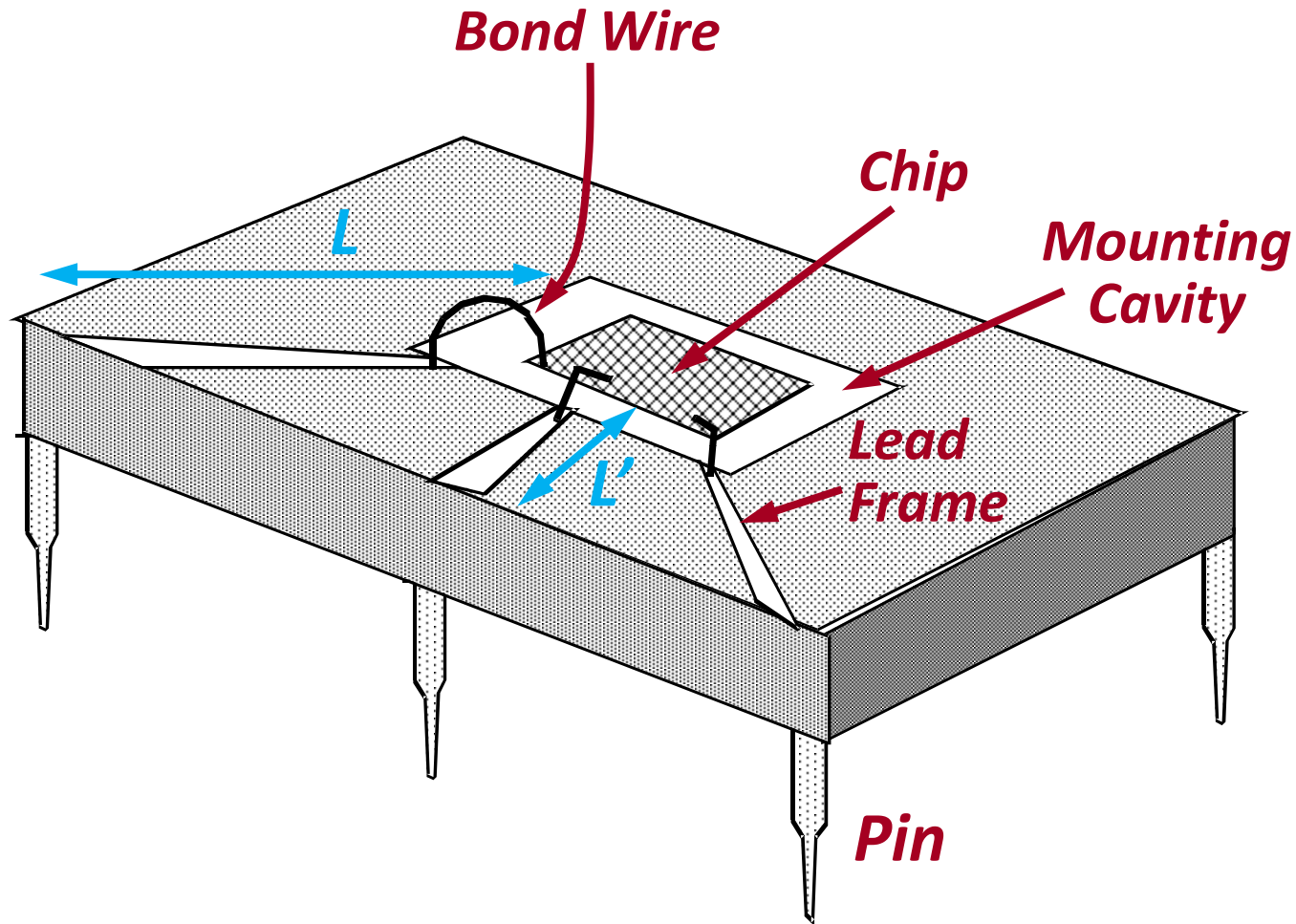
D. Rozgic et al., TBCAS 2019.



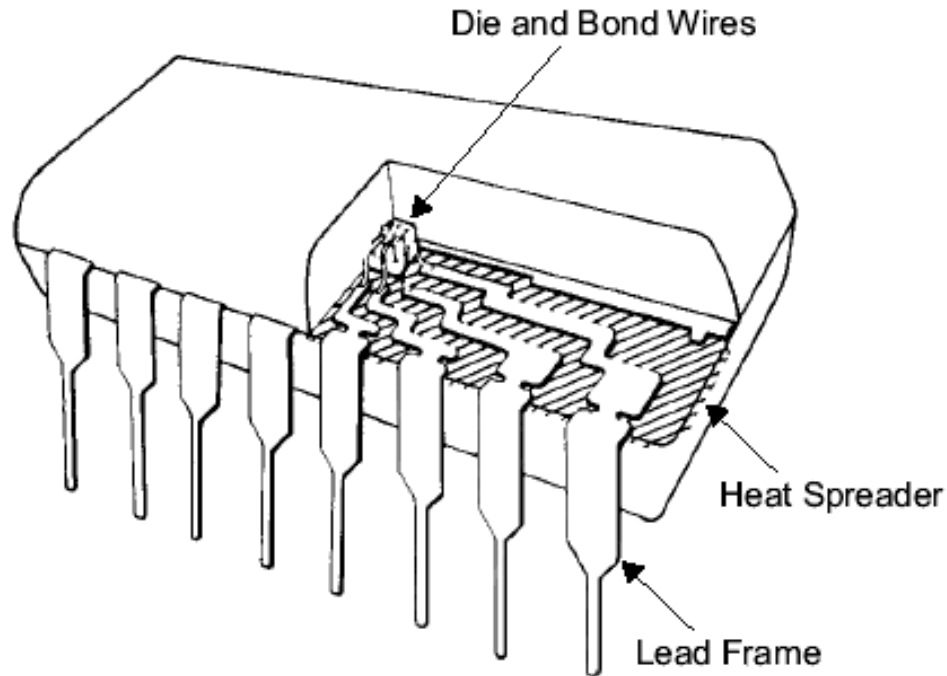
2-layer Si-IF  
(~2μ traces / spaces),  
~9.8μm I/O pad pitch

U. Rathore et al., ISSCC 2022.

# Concept of a Typical Chip Package



# Chip-to-Package Connections

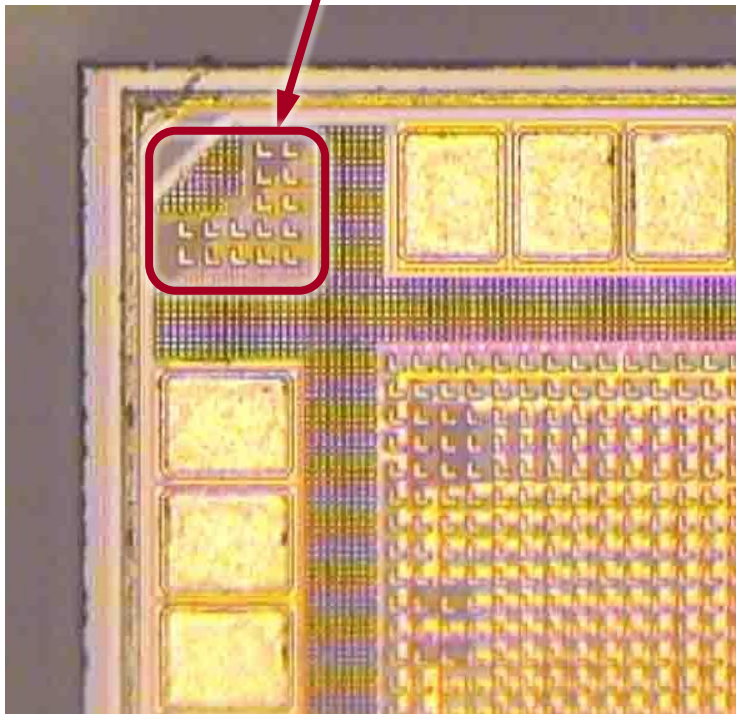


- **Traditionally, chip is surrounded by *pad frame***
  - Metal pads on 50 – 100  $\mu\text{m}$  pitch
  - Gold *bond wires* attach pads to package
  - *Lead frame* distributes signals in package
  - Metal *heat spreader* helps with cooling

# Chip / Package Alignment

Align top-left corners and do simple bonding

Top Left (chip)



Chip

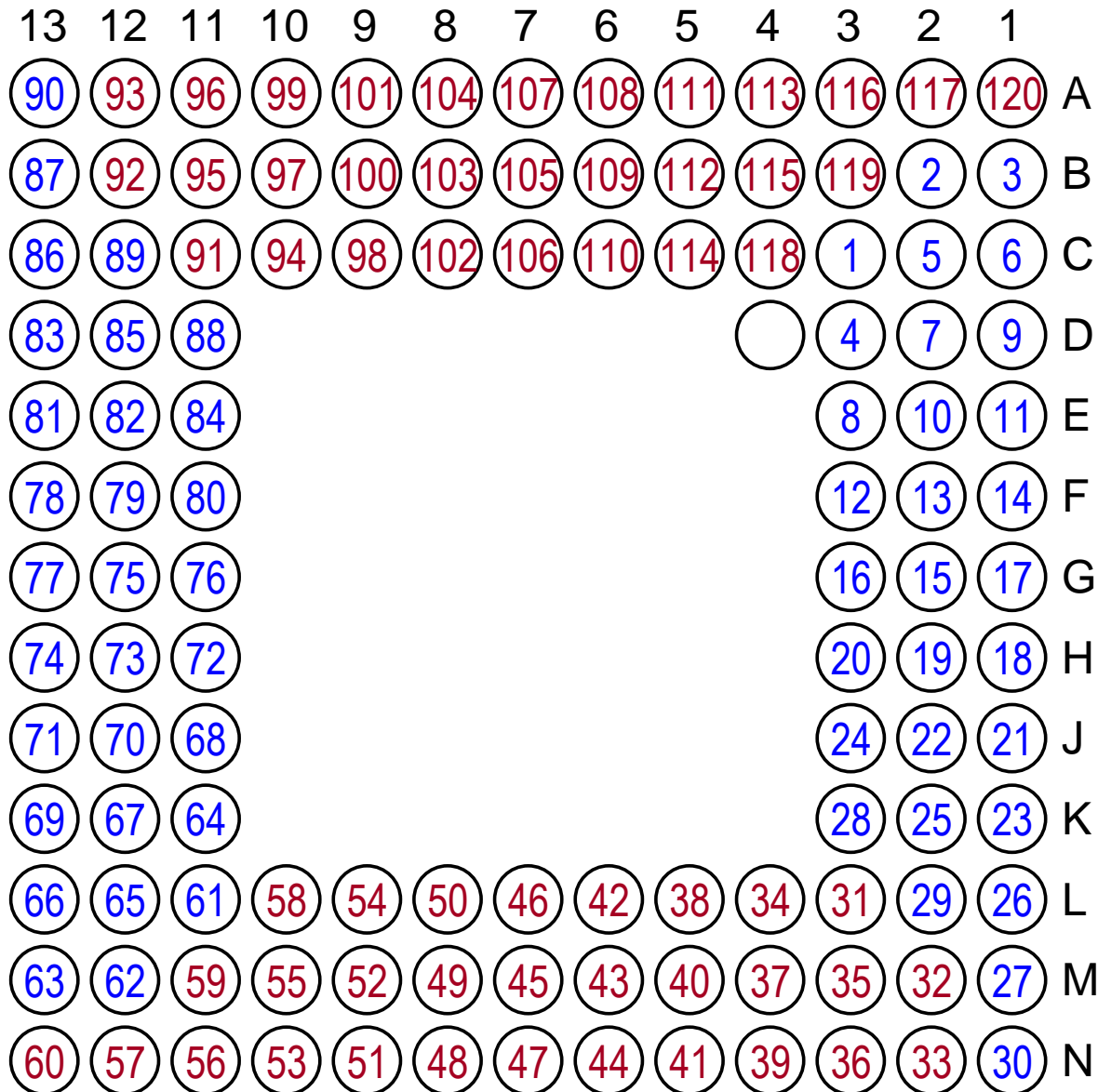
Top Left (index marker)



Package



# Know Pin Locations (PCB Routing)



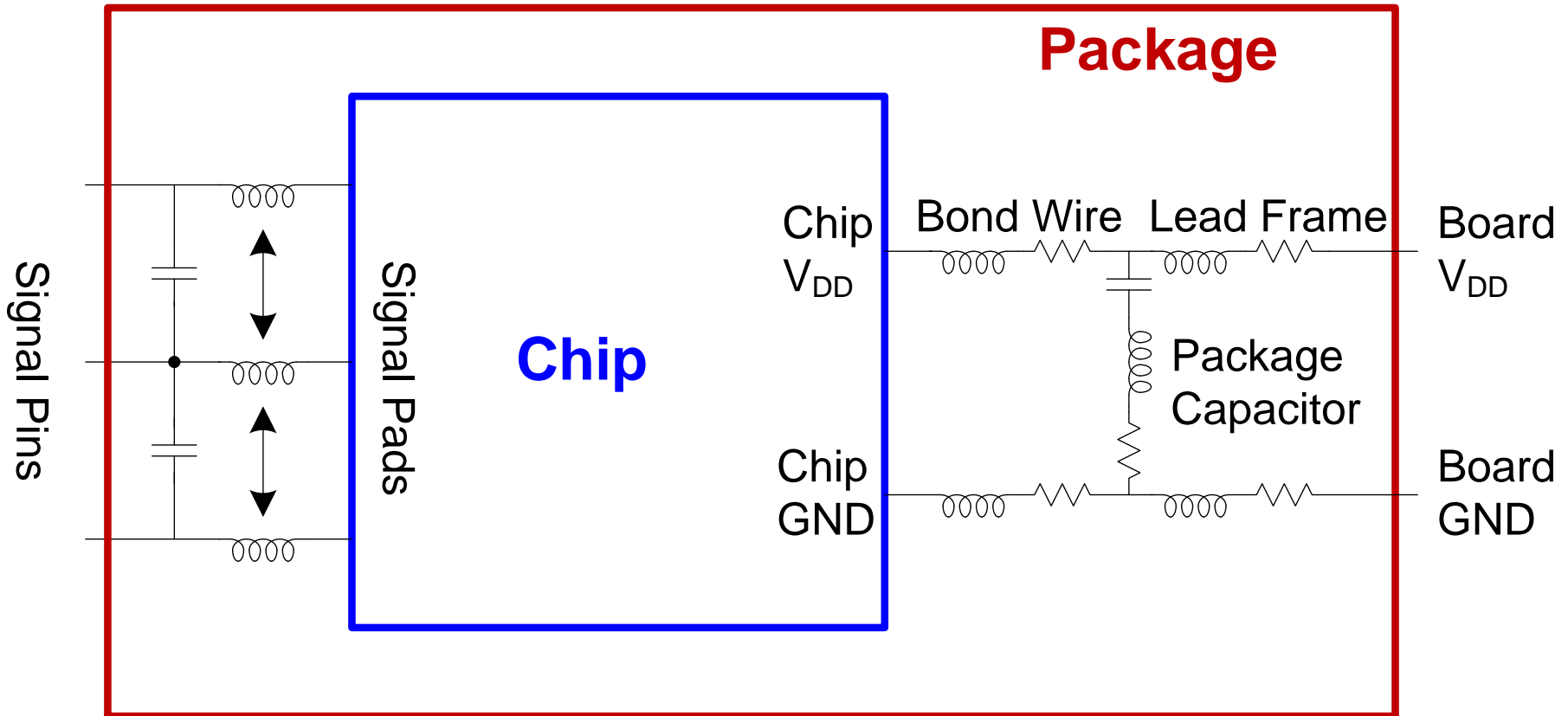
- **120 pins**
  - chip: 3.5mm<sup>2</sup>
- **Ceramic PGA**
  - SSM P/N  
CPG12028  
(~ \$30)
- **Test socket**
  - YAMAICHI  
NP89-19601  
-KS11730  
(~ \$70)

# Advanced Packages

---

- **Bond wires contribute parasitic inductance**
- **Fancy packages have many signal, power layers**
  - Like tiny printed circuit boards
- ***Flip-chip* places connections across surface of die rather than around periphery**
  - Top-level metal pads covered with solder balls
  - Chip flips upside down
  - Carefully aligned to package (done blind!)
  - Heated to melt balls
  - Also called C4

# Package Parasitics



- Use many  $V_{DD}$ , Gnd in parallel
  - Inductance,  $I_{DD}$



# Heat Dissipation: Comparison

---



- 60 W
- Surface area  $\sim 120 \text{ cm}^2$   
(too hot to touch)

- 130 W
- Die area  $\sim 4 \text{ cm}^2$   
(60x higher power density)

# Heat Dissipation

---

- **The heat flows from the transistor junctions through the substrate and package**
  - Can be spread across a heat sink,
  - Then carried away through the air by convection
  - Liquid cooling used in extreme cases (\$\$\$)

## Analogy:

- **Current flow:**  $\Delta V / R$
- **Heat flow:**  $\Delta T / R_{\text{thermal}}$

# Thermal Resistance

---

- $\Delta T = \Theta_{ja} \cdot P$ 
  - $\Delta T$  : temperature rise on chip
  - $\Theta_{ja}$  : thermal R of chip junction to ambient
  - $P$  : power dissipation on chip
- **Thermal resistances combine like resistors**
  - Series and parallel
- $\Theta_{ja} = \Theta_{jp} + \Theta_{pa}$ 
  - Series combination

# Thermal Impedance

---

- **Ceramic pin-grid arrays – 15 to 30 °C/Watt**
- **Plastic Quad Flat Packs – 40 to 50 °C/Watt**
- **Heat dissipation:**
  - Finned heat sinks
  - Embedded metal slugs
- **High-cost packages:**
  - Forced air or liquid cooling through package ducts
  - **Example:** IBM Thermal Conduction Module

# Example 17.1: Thermal Resistance & Power

- Your chip has a heat sink with a thermal resistance to the package of  $4.0^\circ \text{C/W}$   $\left. \vphantom{\Theta_{ja}} \right\} \Theta_{pa}$
- The resistance from chip to package is  $1^\circ \text{C/W}$   $\left. \vphantom{\Theta_{ja}} \right\} \Theta_{jp}$
- The system box ambient temp. may reach  $55^\circ \text{C}$
- The chip temperature must not exceed  $100^\circ \text{C}$

$$\Theta_{ja} = 5^\circ \text{C/W}$$

$$\Delta T = 45^\circ \text{C}$$

$$\Delta T = \Theta_{ja} \cdot P_{max}$$

- What is the maximum chip power dissipation?

- $\Theta_{ja} = \Theta_{jp} + \Theta_{pa} = 1 + 4 = 5^\circ \text{C/W}$

- $\Delta T = 100 - 55 = 45^\circ \text{C}$

- $P = \Delta T / \Theta_{ja} = 45/5 = 9 \text{ W}$

# Power Distribution Network Functions

---

- Carry current from pads to transistors on chip
- Maintain stable voltage with low noise
- Provide average and peak power demands
- Provide current return paths for signals
- Avoid electromigration & self-heating wearout
- Consume little chip area and wire
- Easy to lay out

# Power Supply Drop/Noise

---

$$R_{\text{power}} \cdot I_{\text{VDD}} = \Delta V_{\text{DD}}$$

**Supply noise** is variations in power supply that manifest as noise onto the logic gates

- Power supply **wiring resistance** creates voltage variations with current surges
- The **current surge** for static CMOS depend on dynamic behavior of circuit

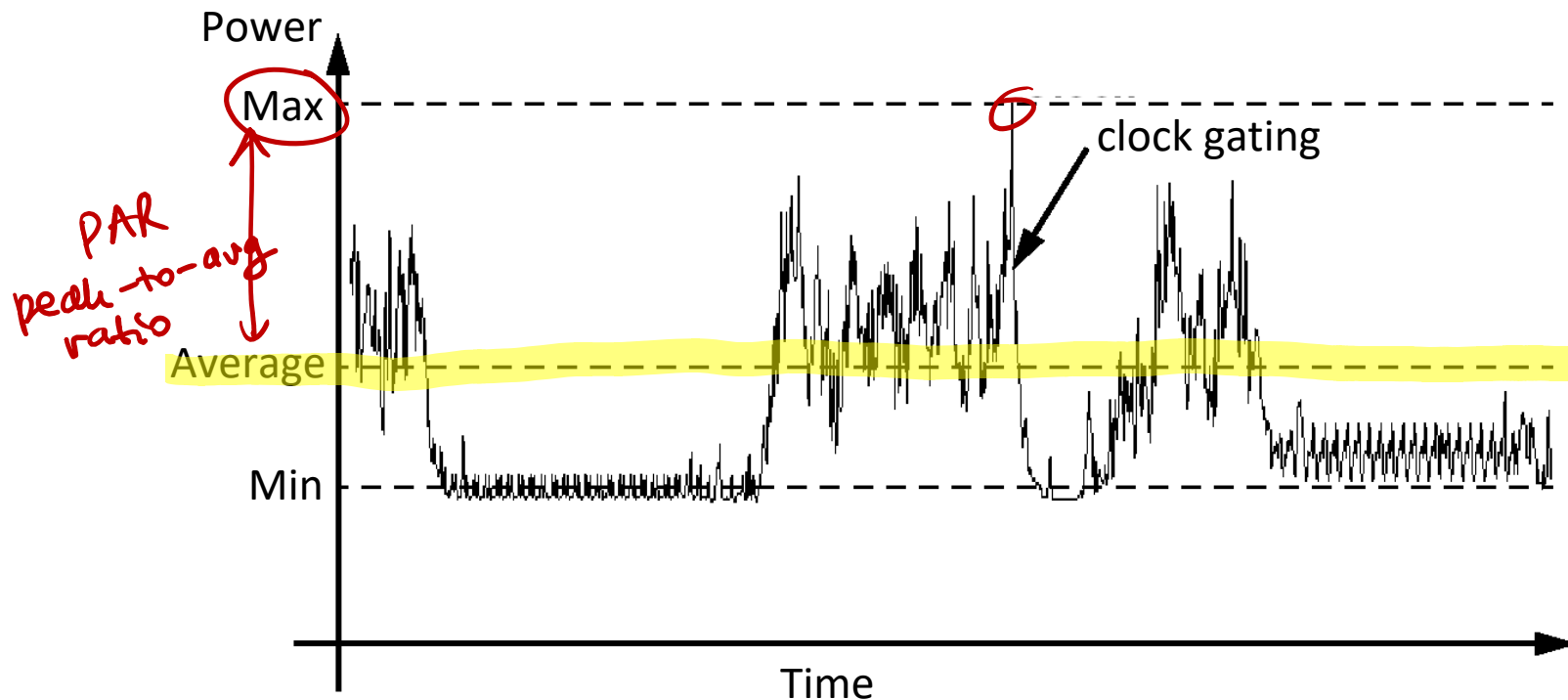
# Tackling the $V_{DD}$ Drop

---

- 1** •  $V_{DD}$ -Gnd capacitance
  - Based on total max  $C_{switched}$  (10x)  
*& provides  $I_{peak}$*
- 2** • Redesign power/ground network to reduce resistance
  - Based on max  $I_{DD}$  required by each block  
*↑ lowers  $R_{pwr}$*
- 3** • Adjust activity to another clock cycle to reduce peak current
  - Scheduling



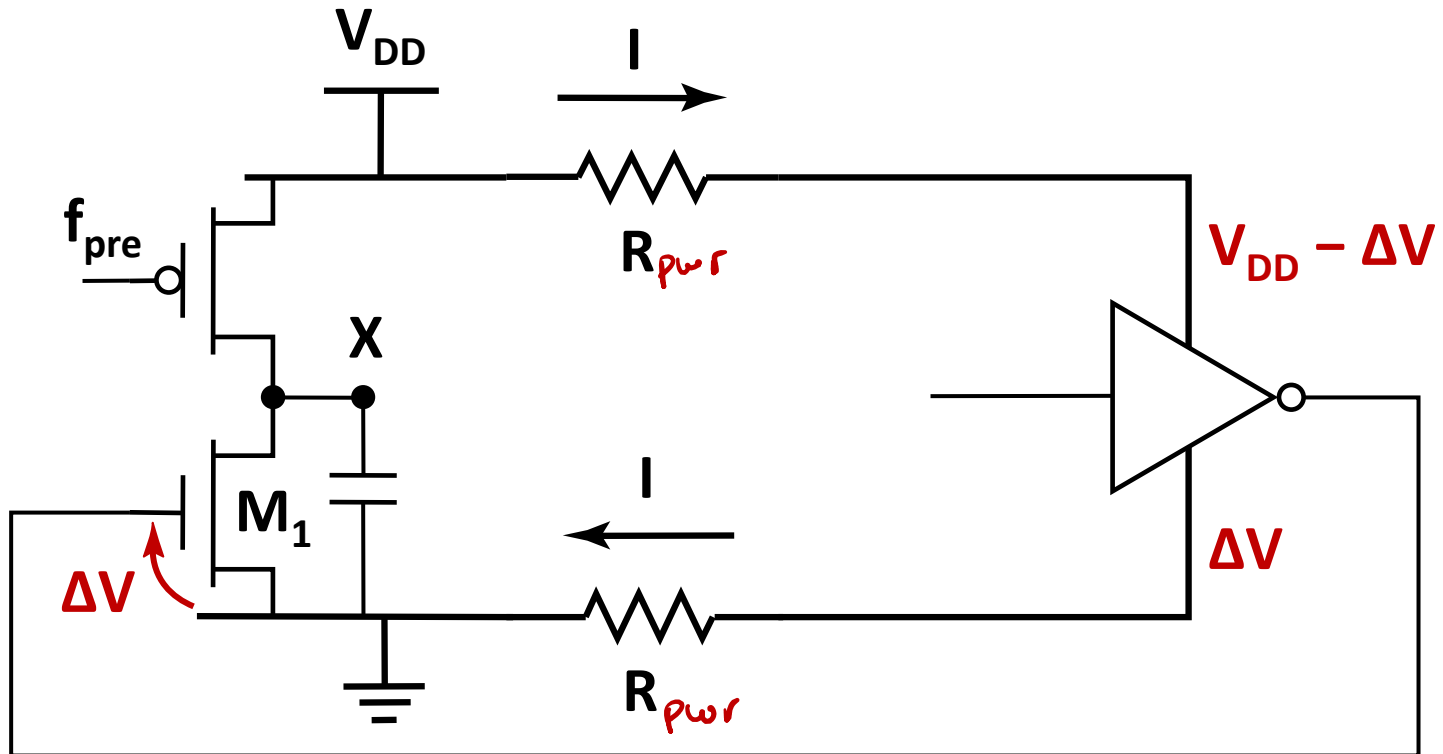
# Power Requirements



- $V_{DD} = V_{DDnominal} - V_{droop}$
- Want  $V_{droop} < \pm 10\%$  of  $V_{DD}$
- $I_{DD}$  changes on many time scales

- Sources of  $V_{droop}$ 
  - IR drops ✓
  - L di/dt noise ✓

# Issue #1: RI Introduced Noise

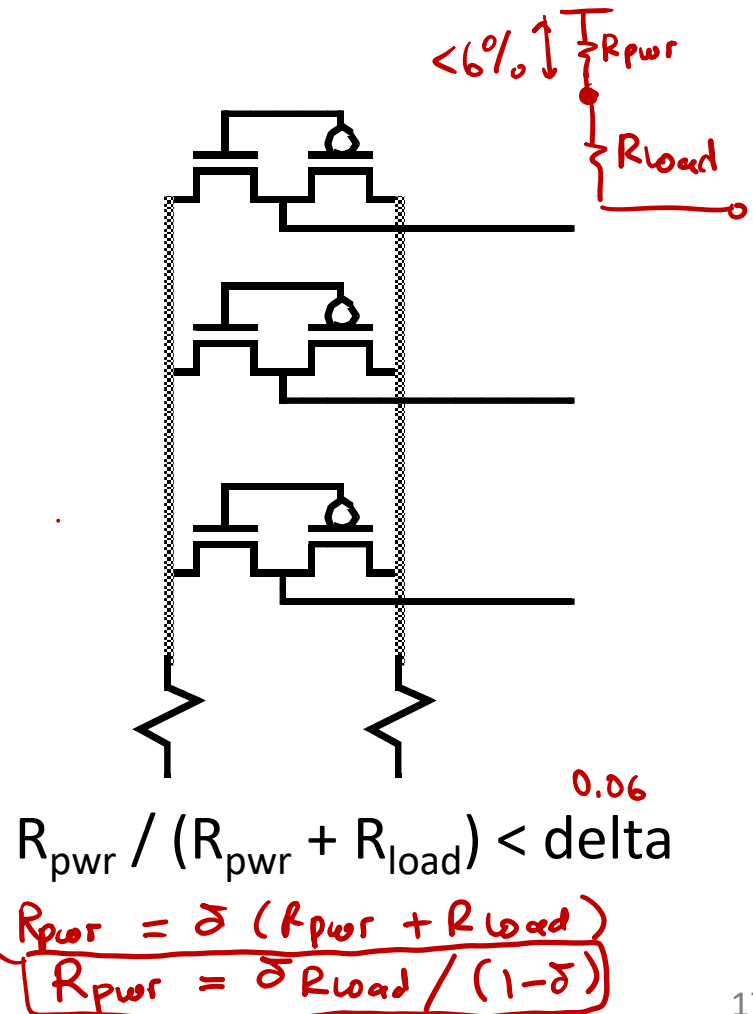


# Example 17.2: Power IR Drop

- Drive a 32-bit bus, total load of each wire: 2pF,  $R = 0.125\Omega/\text{square}$ , want delay  $\sim 0.5\text{ns}$ ,  $< 6\%$  drop

- R for each transistor needs to be  $< 0.25\text{ k}\Omega = R$ 
  - To meet  $RC = 0.5\text{ ns}$   $\leftarrow 2\text{ pF}$
- Effective R of bits together is  $250/32 = 7.5\ \Omega = R_{\text{load}}$

- For  $< 6\%$  drop, Power R must be  $< 0.48\ \Omega$ 
  - That is only 4 squares



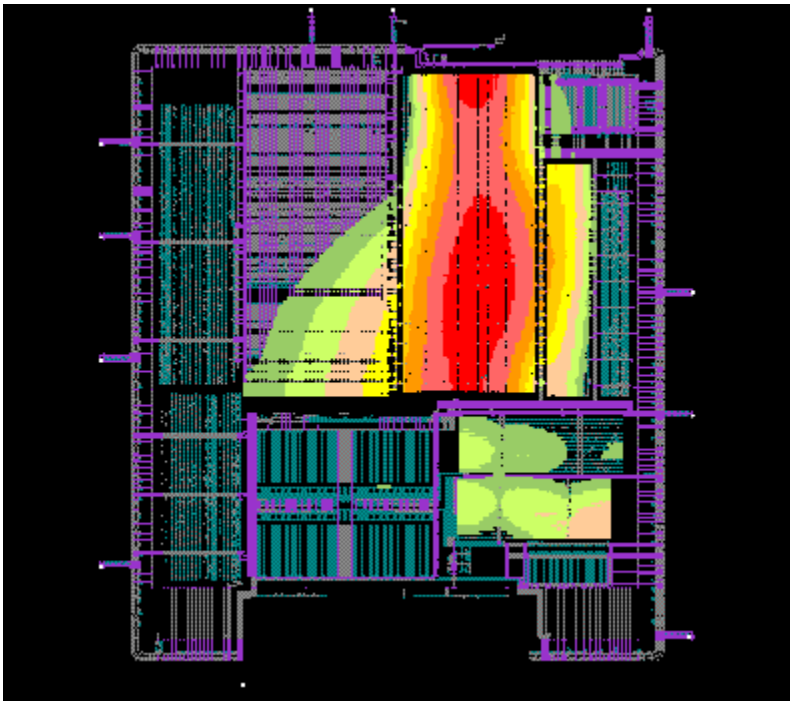
# Must Support Total Power

---

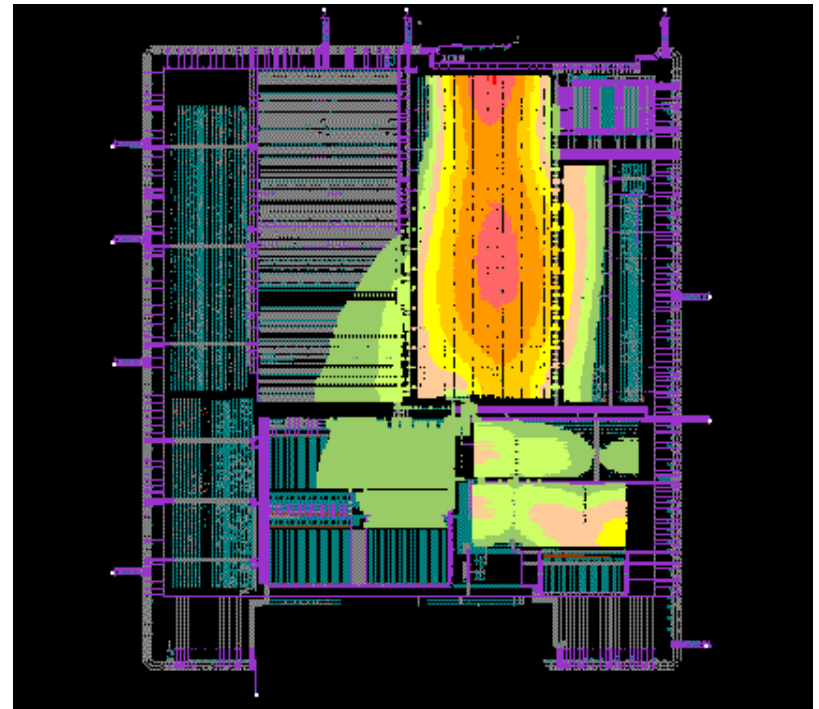
- **Chips today dissipate 5-50W**
- **Implies total current is 5-50A (Power = IV)**
  - Supply is now ~1V!
- **Very big problem currently**
- **Use many supply pins (@ few mA each),  
and wide wires for low R**
- **Grids of high-level metal for power is a must!**
  - Thicker metal... lower R

# Resistance and Power Distribution Problem

Before



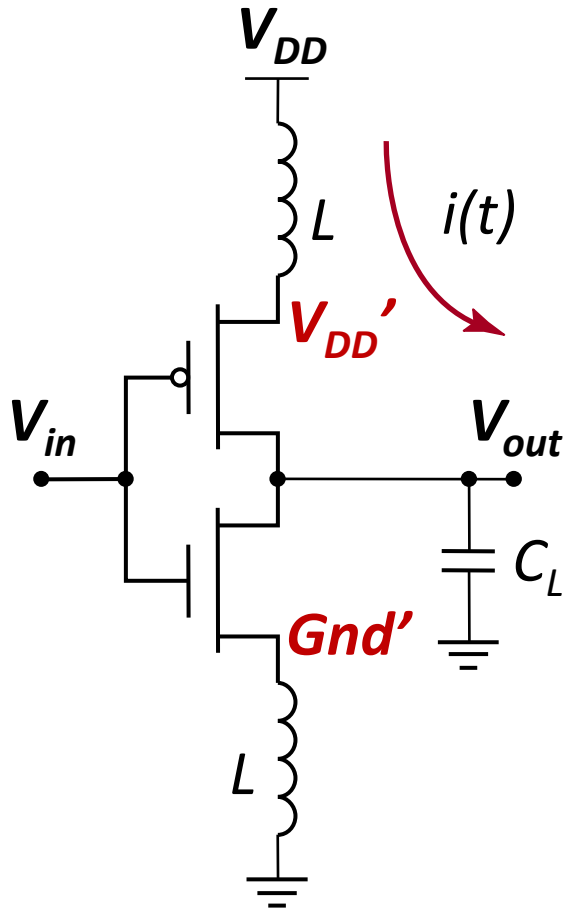
After



- Requires fast and accurate peak current prediction
- Heavily influenced by packaging technology

*Source:  
Cadence*

# Issue #2: $L \, di/dt$



## Impact of inductance on supply:

- Change in current induces the change in voltage
- Longer supply lines have larger  $L$

## Example 17.3: L di/dt Calculation

- (12.3.3 W&H) A 1GHz chip transitions from idle (20 A) to full power (60 A) operation in a single cycle
- If the power supply has 20 pH of series inductance, estimate the power supply noise caused by this transition

### Solution:

- $\Delta I/\Delta t = (60\text{A} - 20\text{A})/1\text{ns} = 40 \text{ GA/s}$
- The inductive noise is:  $L \Delta I/\Delta t = 0.8 \text{ V}$ 
  - Unacceptable in a low-voltage process
  - The chip needs internal bypass capacitance

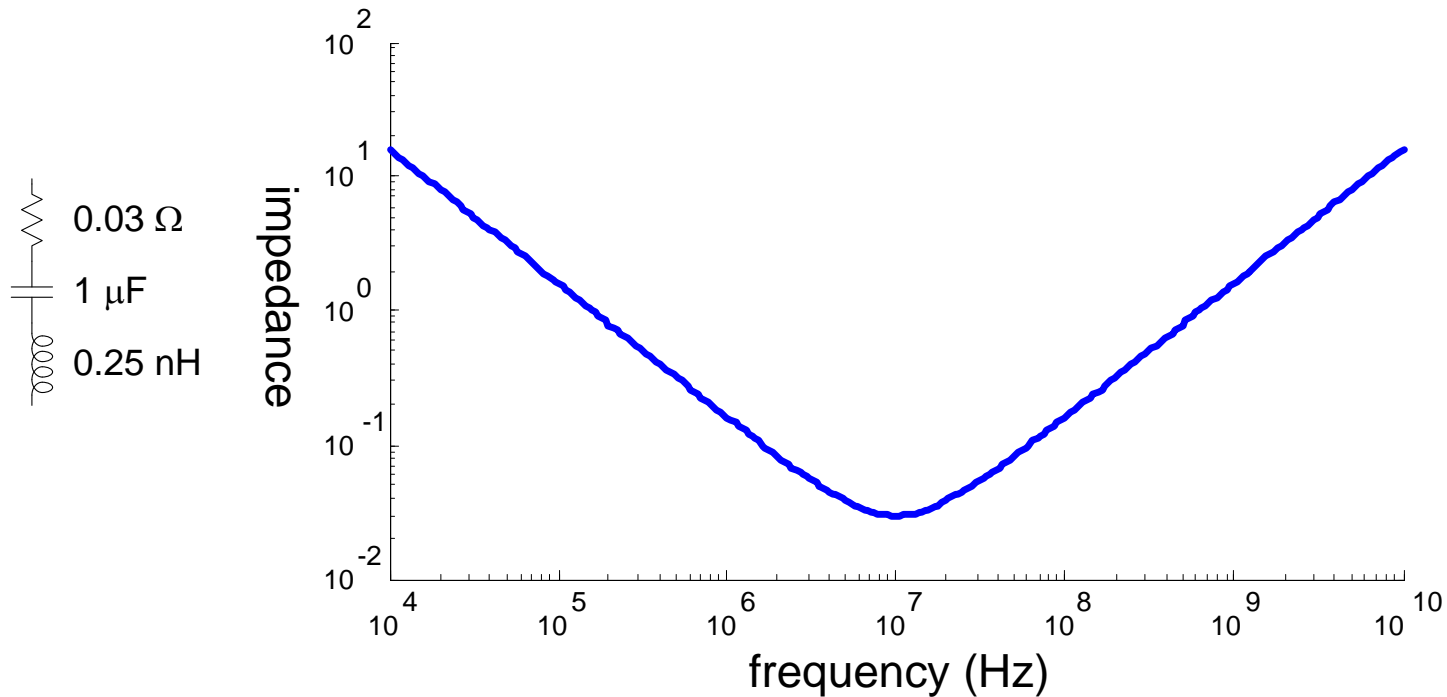
# Design Techniques to Address $L di/dt$

---

- Separate power pins for I/O pads and chip core
- Multiple power and ground pins
- Position of power and ground pins on package
- Increase  $t_r$  and  $t_f$
- Advanced packaging technologies
- Decoupling capacitances on chip and on board



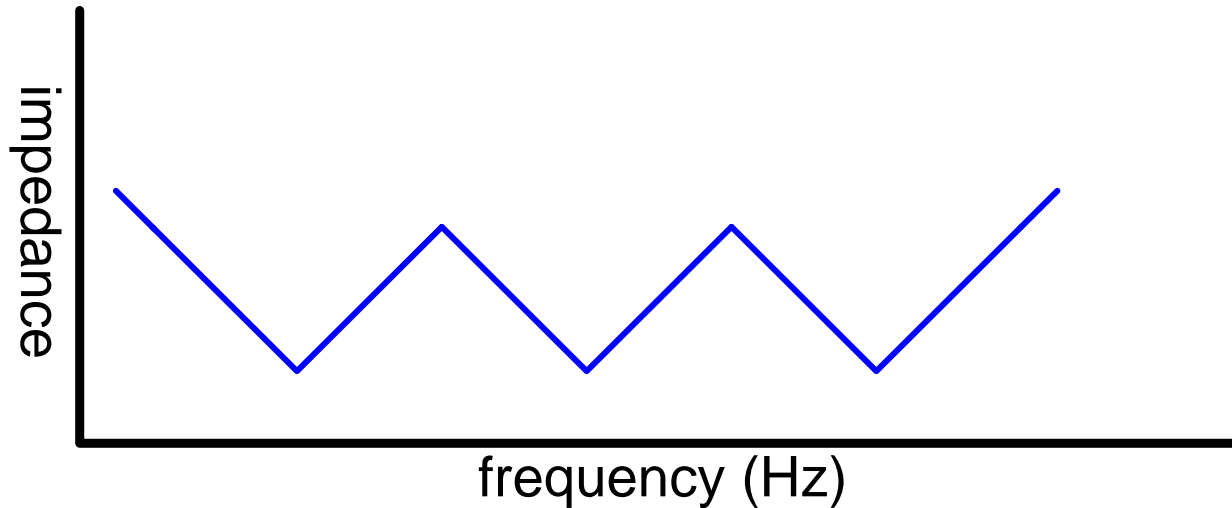
# Bypass Capacitors



- Need low supply impedance at all frequencies
- Ideal caps have impedance decreasing with  $\omega$
- **Real caps have parasitic R and L**
  - Leads to resonant frequency of capacitor

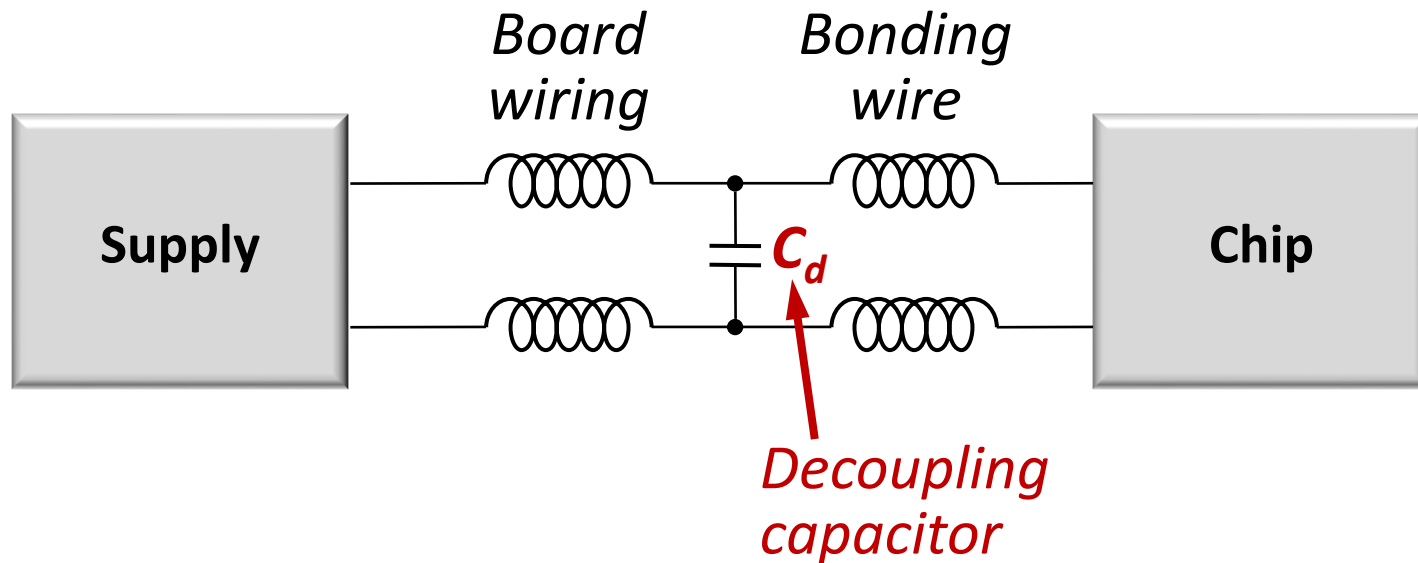
# Use Multiple Capacitors in Parallel

---



- Choose caps to get **low Z at all frequencies**
- Large  $C$  near regulator : low  $Z$  at low frequencies
  - But also has a low self-resonant frequency
- Small capacitors near chip and on chip have low impedance at high frequencies

# Decoupling Capacitors

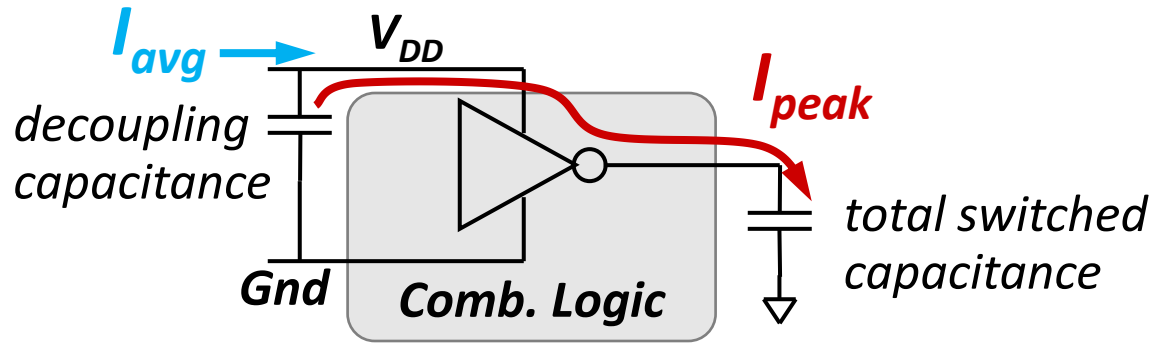


## Decoupling capacitors are added:

- On the board (right under the supply pins)
- On the chip (under the supply straps, near large buffers)

# On-chip Decoupling Capacitance

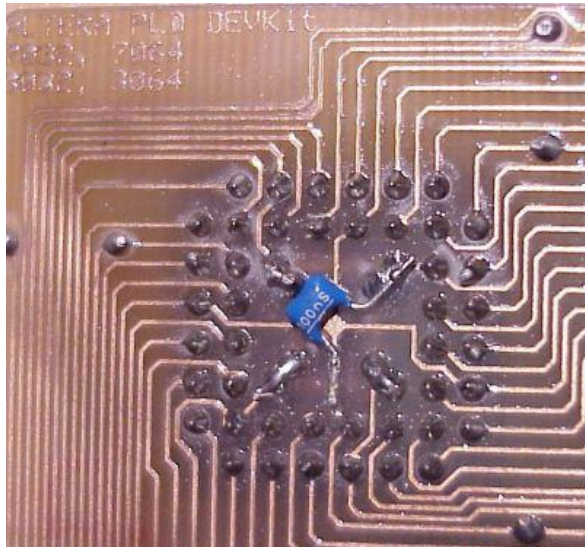
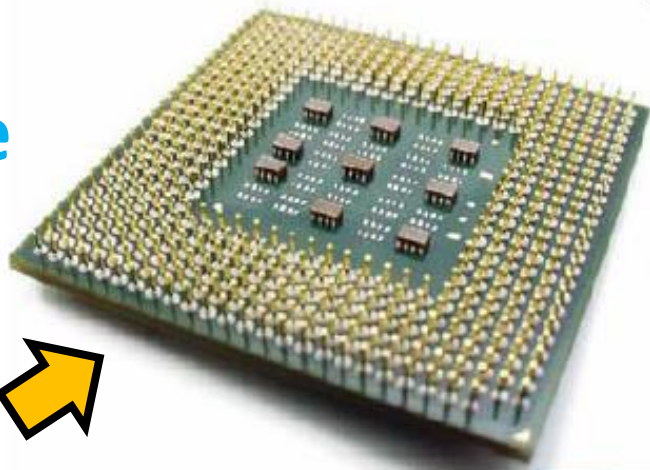
- **Static CMOS logic dynamically switches (no dc current)**
  - Supply just needs to provide the average current
  - Peak current needs to come from nearby capacitance



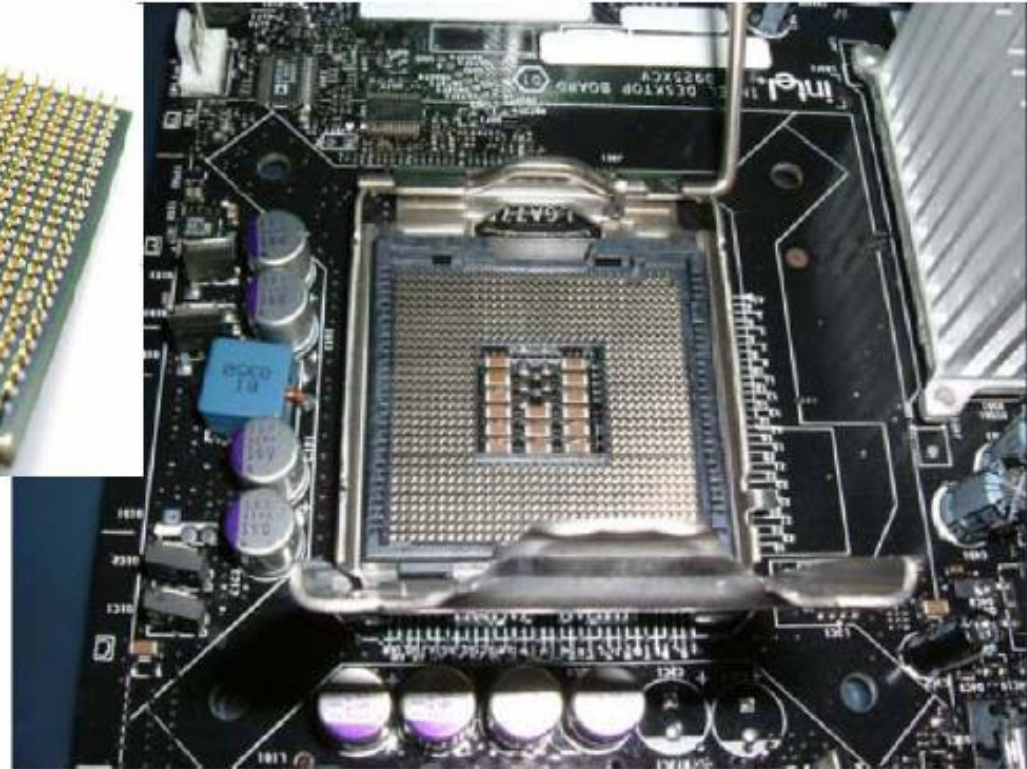
- **Basically the same as charge sharing**
  - Use  $C_{decoup} > 10 C_{switched}$  to guarantee  $< 10\% V_{DD}$  drop
  - Put capacitance near load with little resistance
    - Part of the PnR tool

# Bypass Capacitances in Real Life

Package  
bypass



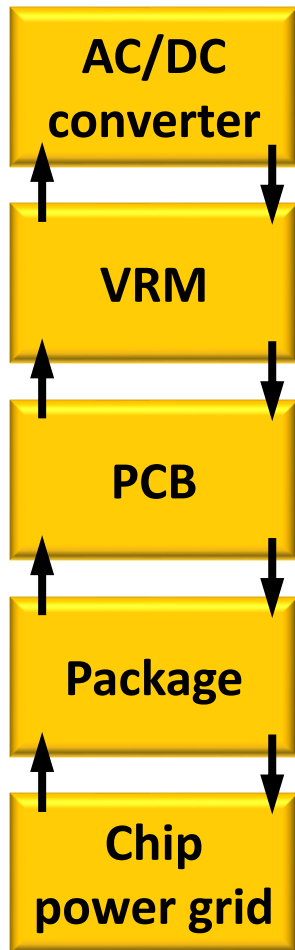
100 nF decap  
FPGA chip



PCB bypass

Image courtesy of  
M. Horowitz & K. Mai

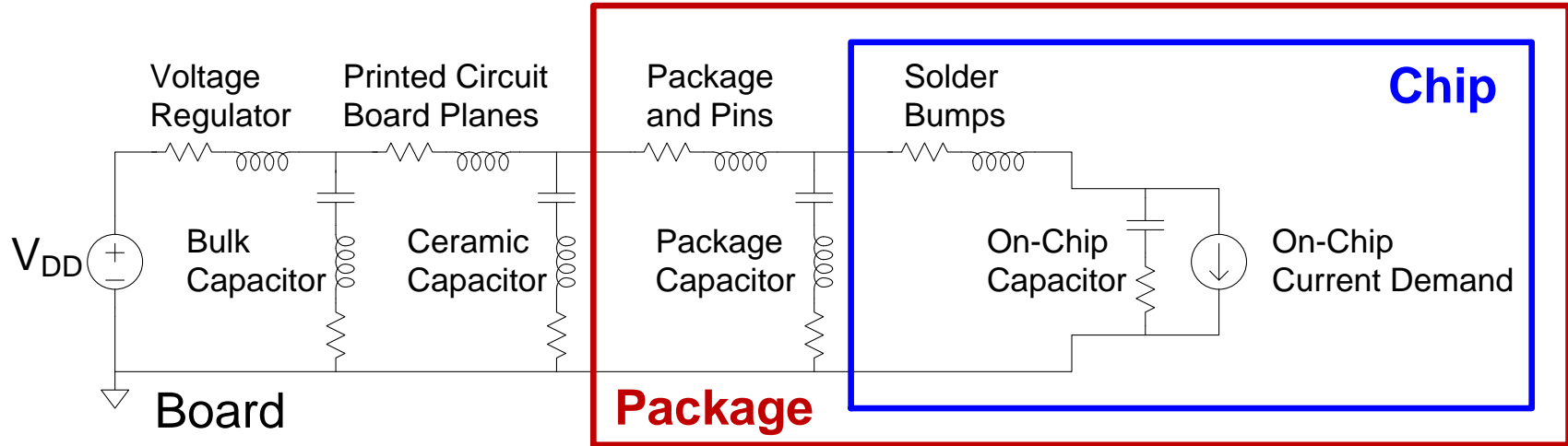
# Power Distribution Network



- **AC/DC converter**
  - Usually 110 VAC to 12 or 5 VDC in PCs
- **Voltage Regulator Module**
  - Converts one DC level to another (5V to 1V)
- **Printed Circuit Board**
  - Planes send current from VRM to the package
  - Planes have caps for bypass; use discretes too
- **Package**
  - Deliver current to the chip using balls or bonds
  - Can use bypass caps on the package as well
- **Chip power grid**
  - Use device bypass capacitors

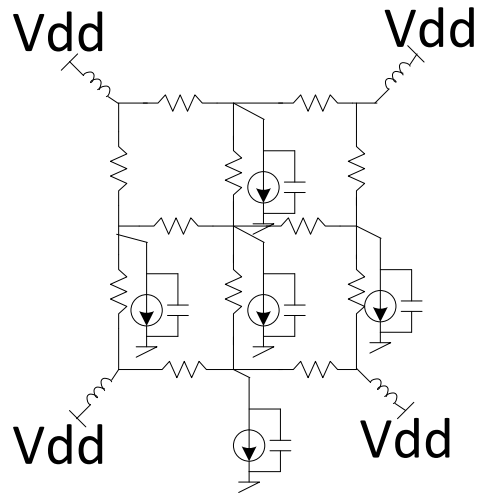
Courtesy:  
M. Horowitz

# Power System: Lumped Model

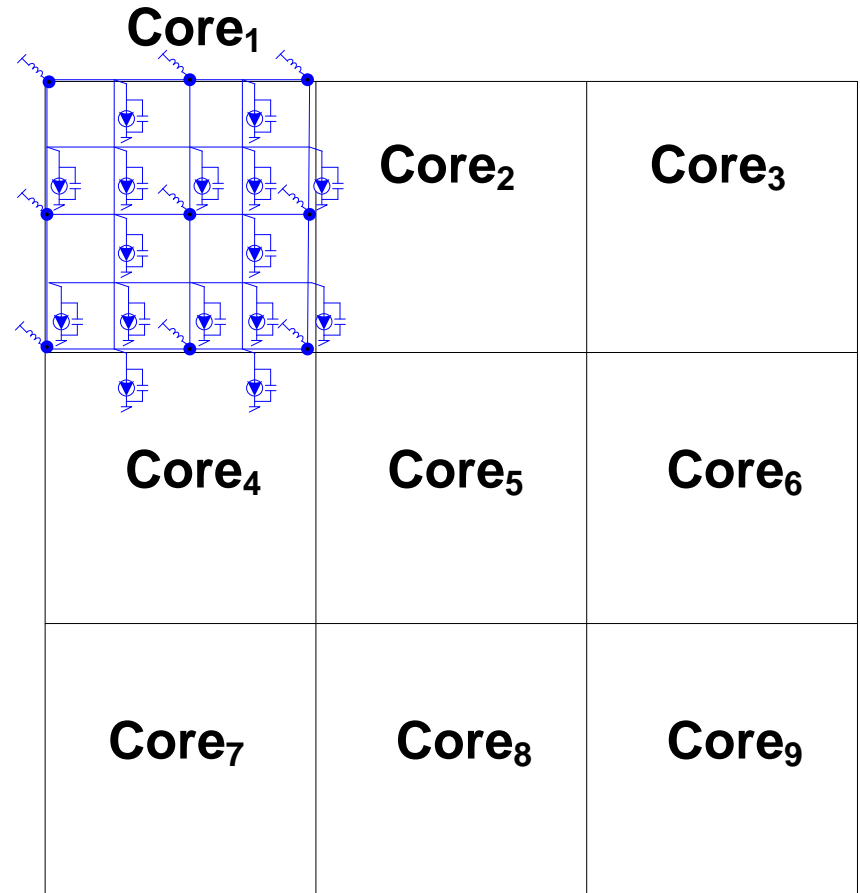


- **Power comes from regulator on system board**
  - Board and package add parasitic  $R$  and  $L$
  - Bypass capacitors help stabilize supply voltage
  - But capacitors also have parasitic  $R$  and  $L$
- **Simulate system for time and frequency responses**

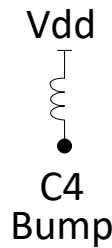
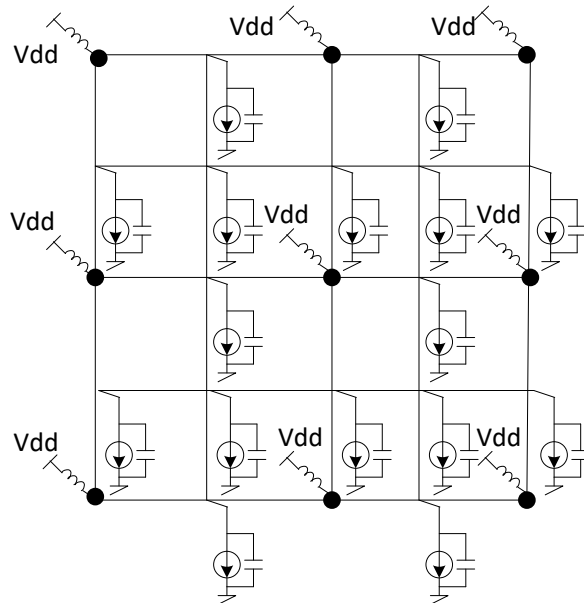
# On-chip Power Network: Distributed Model



**Base grid**



**Global grid**



**Core grid**

Courtesy: A. Todri

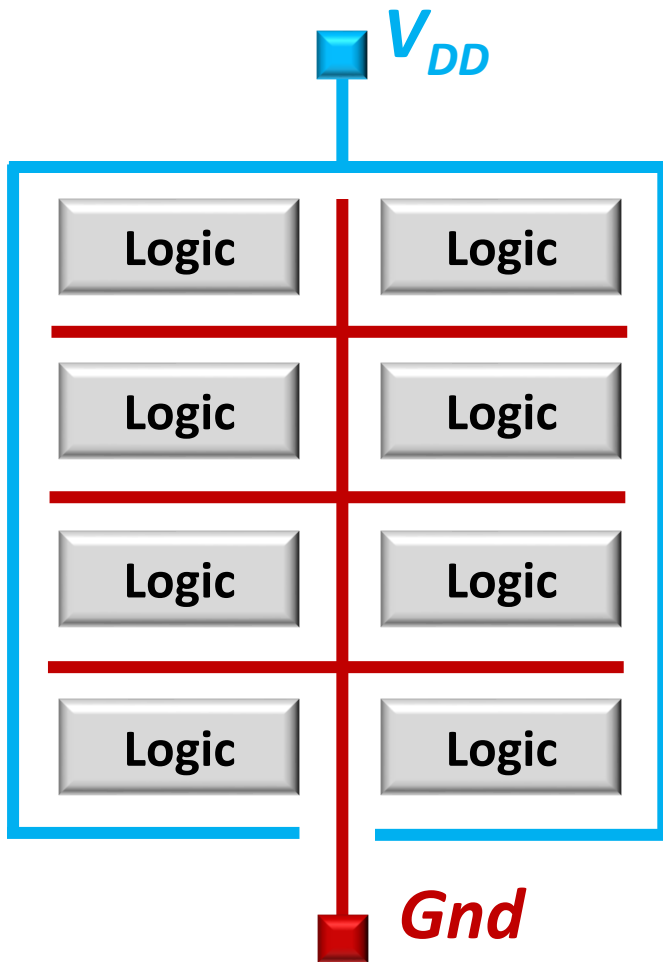


# Power Distribution Strategy

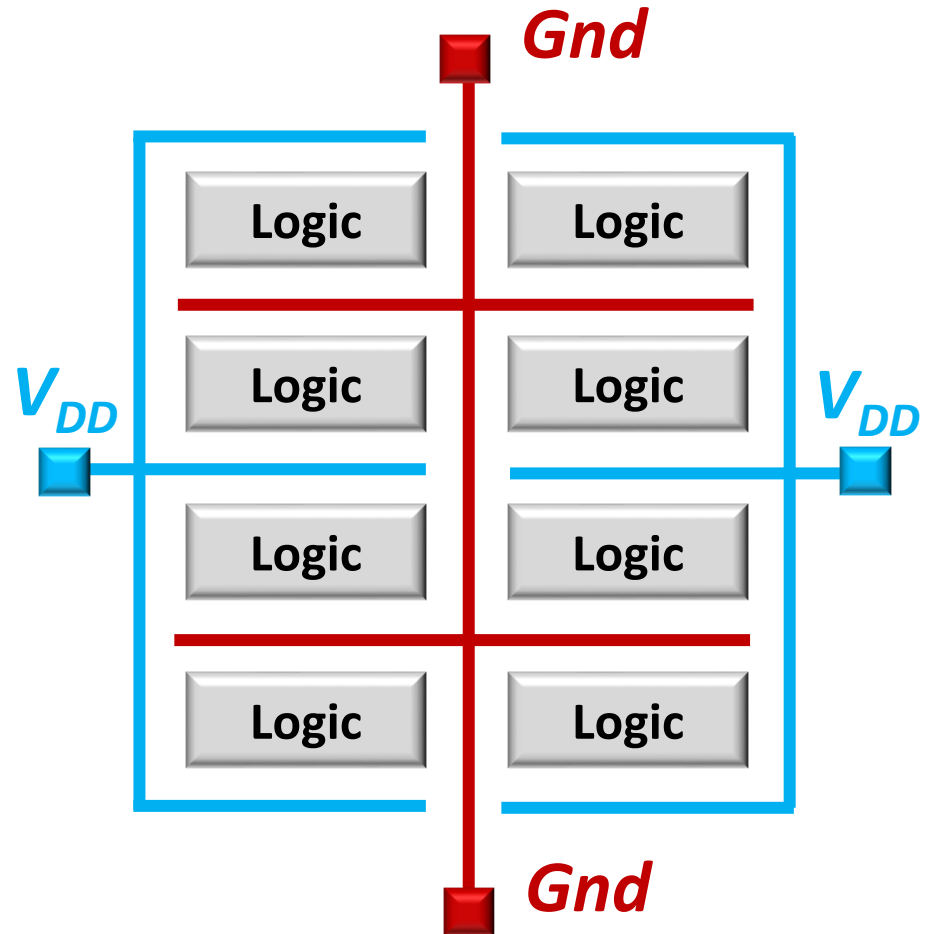
---

- **Low-level distribution is in Metal 1**
- **Power has to be “strapped” in higher layers of metal**
- **The spacing is set by IR drop, electromigration, inductive effects**
- **Always use multiple contacts on straps**

# Power and Ground Distribution



- Finger-shaped network



- Multiple supply pins

# Power and Ground Need Wider Wires

---

**Power needs to be distributed to all the cells in the circuit**

- **A tree**
  - Trunk of the tree must supply current to all branches
- **R in these lines must be very small, since when a gate switches, its current flows through the supply lines**
  - If R of the supply lines is too large, the voltage supplied to gates will drop, which can cause the gate to malfunction
  - Usually you don't want the supply to change more than 5-10% due to supply resistance

# Power Must be on the Metal Layer

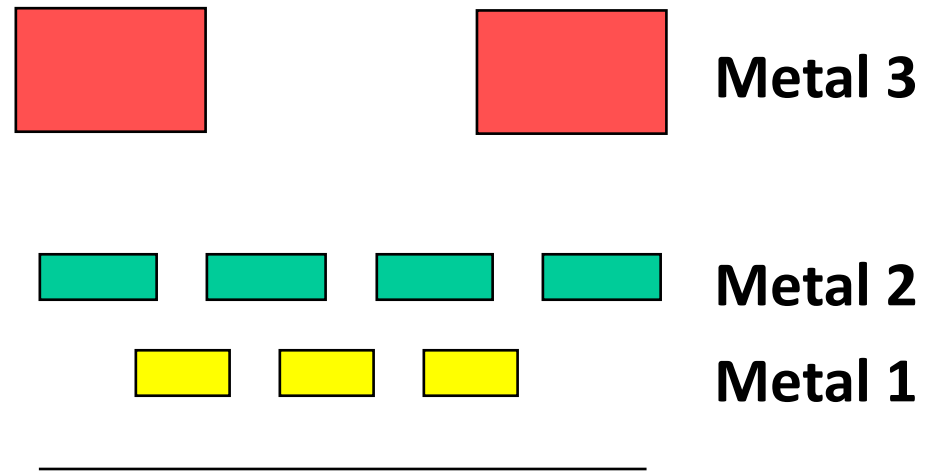
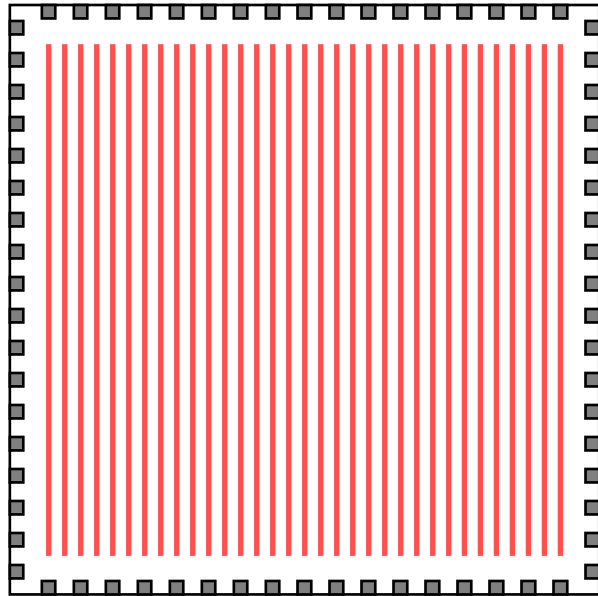
---

## Is that enough?

Usually, they have to be wider too

- $R_{\text{trans}}$  is much greater (by  $10^5$ ) than  $R_{\text{metal}}$
- But one builds wide devices, and long wires
- And in a chip there are many devices connected in parallel to the supplies
  - So you need still need to be careful even with metal layers, and make the special wires wide enough

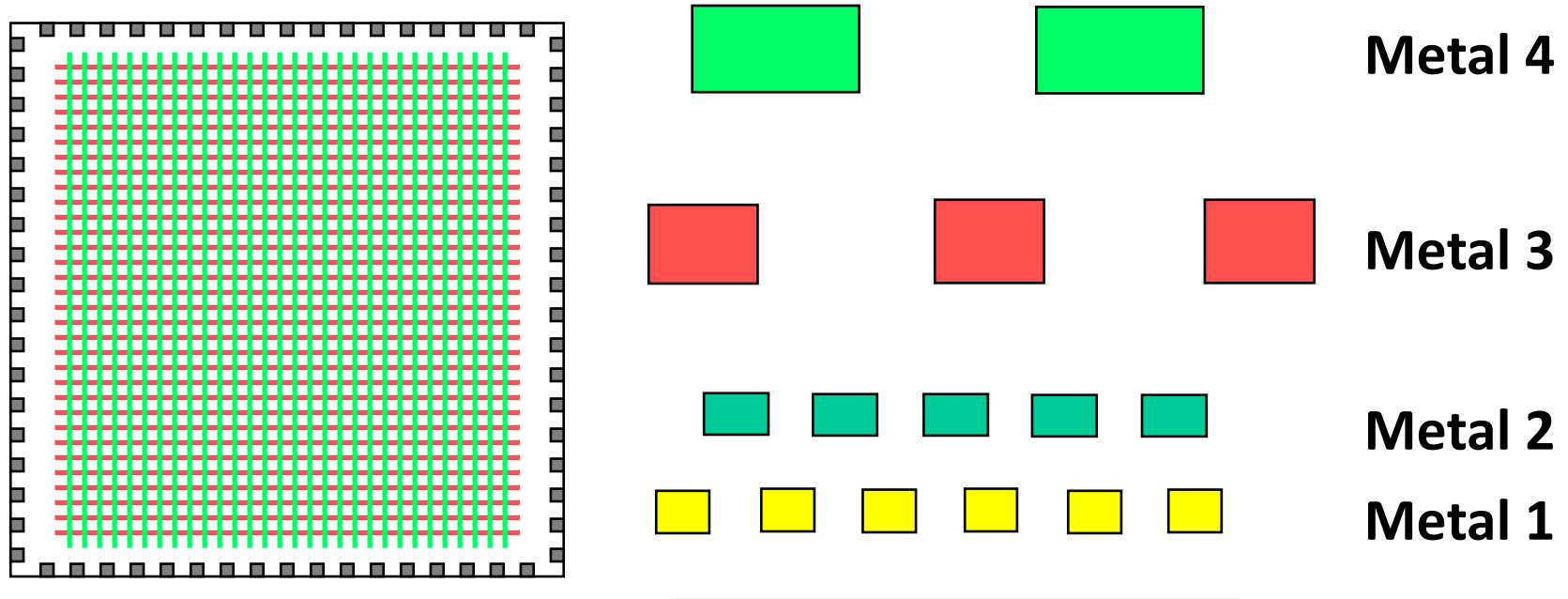
# 3-Metal Layer Approach (EV4)



- Power supplied from **two sides** in M3
- M2 used to form power grid
- **90% of M3 used for power/clk routing**

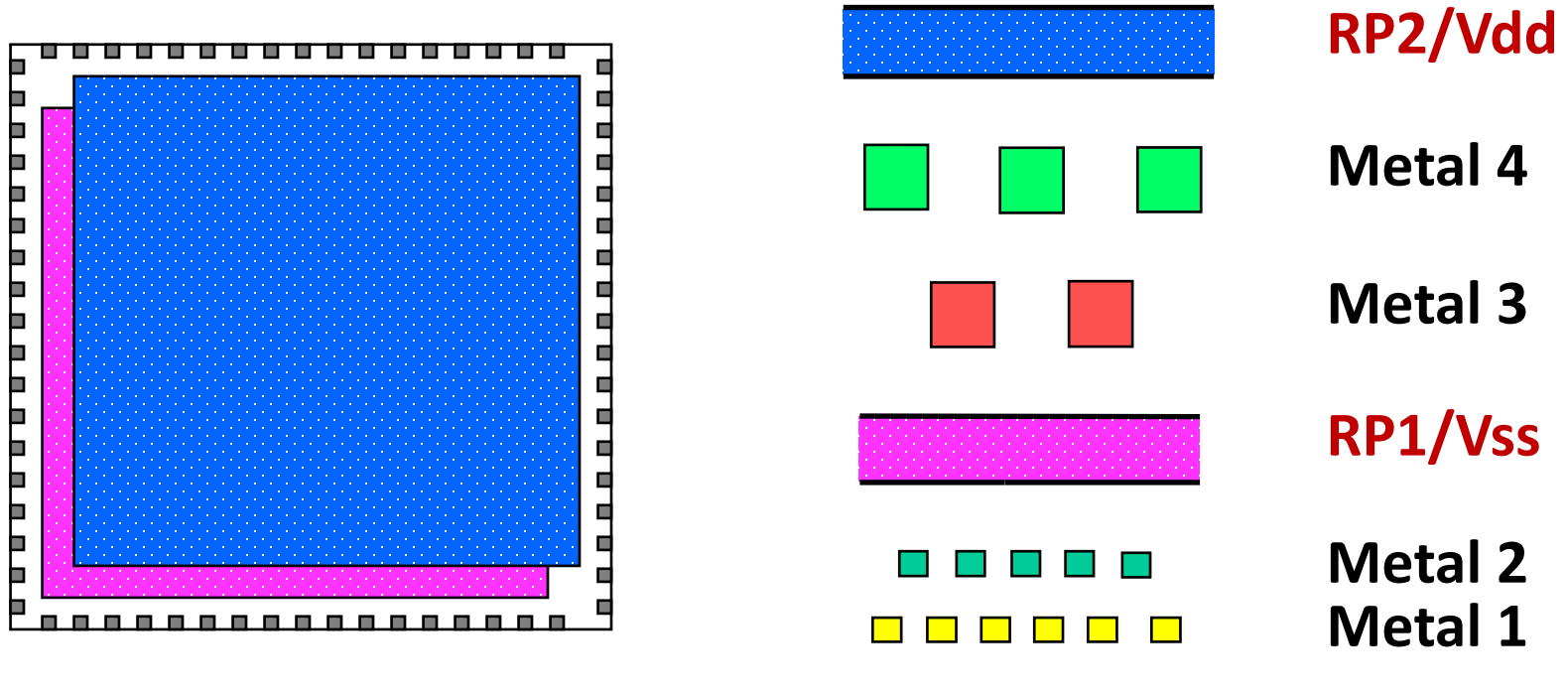
*Courtesy: Compaq*

# 4-Metal Layers Approach (EV5)



- Power supplied from **four sides**
- Grid strapping done all in coarse metal
- **90% of M3/M4 for power/clk routing**

# 6 Metal Layer Approach (EV6)



- **Two reference-plane metal layers**
- **Significantly lowers resistance of grid**
- **Lowers on-chip inductance**

*Courtesy: Compaq*

# Decoupling Capacitor Ratios

Processor	EV4	EV5	EV6
Technology	0.75 $\mu$ m	0.5 $\mu$ m	0.35 $\mu$ m
Year	1992	1995	1998
Clock rate	200 MHz	350 MHz	600 MHz
$C_{\text{switched}}$	12.5 nF	13.9 nF	34 nF
$C_{\text{decoupling}}$	128 nF	160 nF	320 nF

**Decoupling/switching capacitance  $\sim 10x$**

*Source: B. Herrick (Compaq)*



# EV6 De-coupling Capacitance: Example

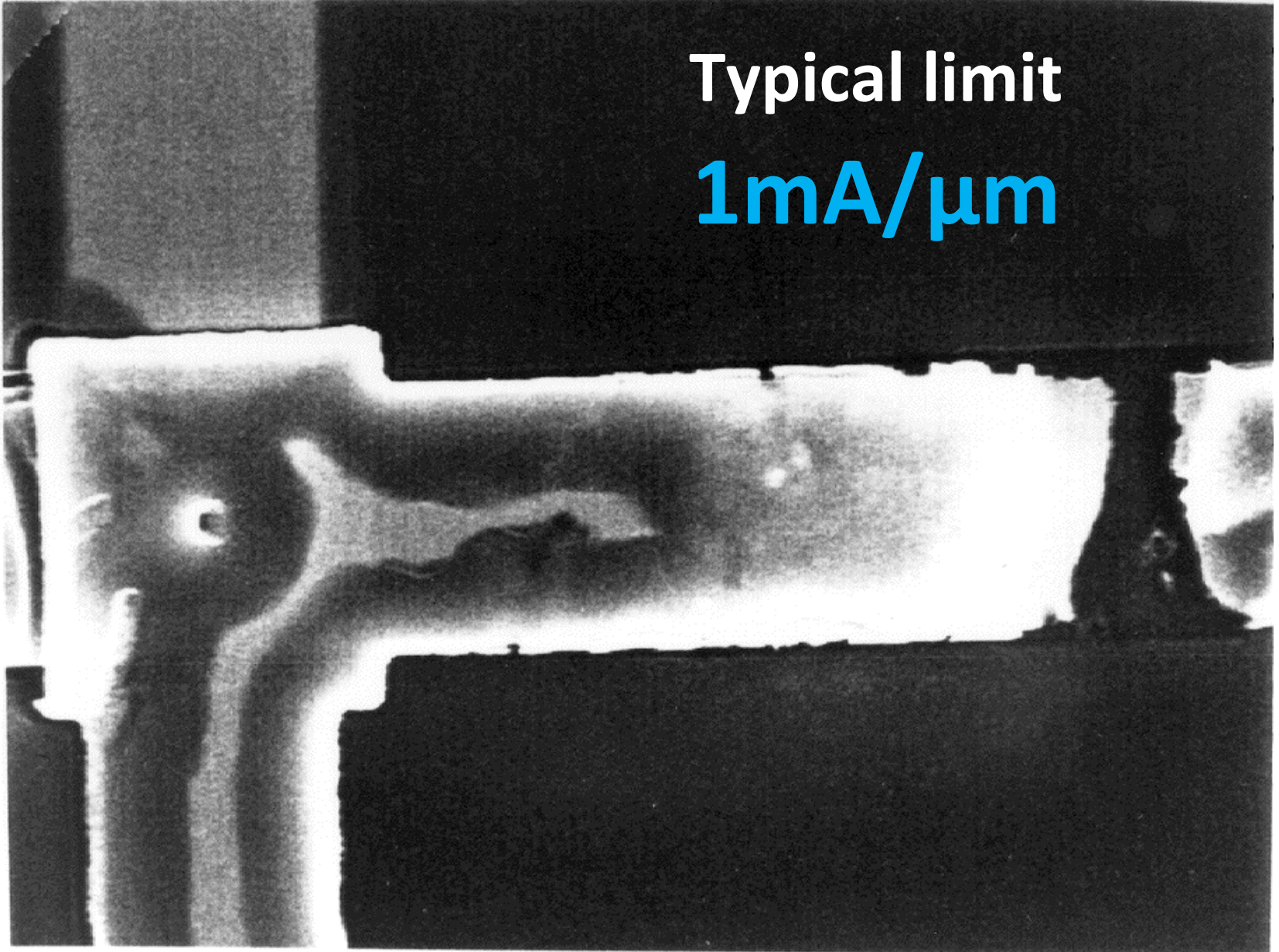
---

Design for  $\Delta I_{DD} = 25 \text{ A @ } V_{DD} = 2.2 \text{ V}, f = 600 \text{ MHz}$

- **0.32- $\mu\text{F}$  of on-chip de-coupling capacitance**
  - Under major busses and around major gridded clock drivers
  - Occupies **15-20% of die area**
- **1- $\mu\text{F}$  2-cm<sup>2</sup> Wirebond Attached Chip Capacitor (WACC) significantly increases “Near-Chip” de-coupling**
  - **160 Vdd/Vss bondwire pairs** on the WACC to **minimize inductance**

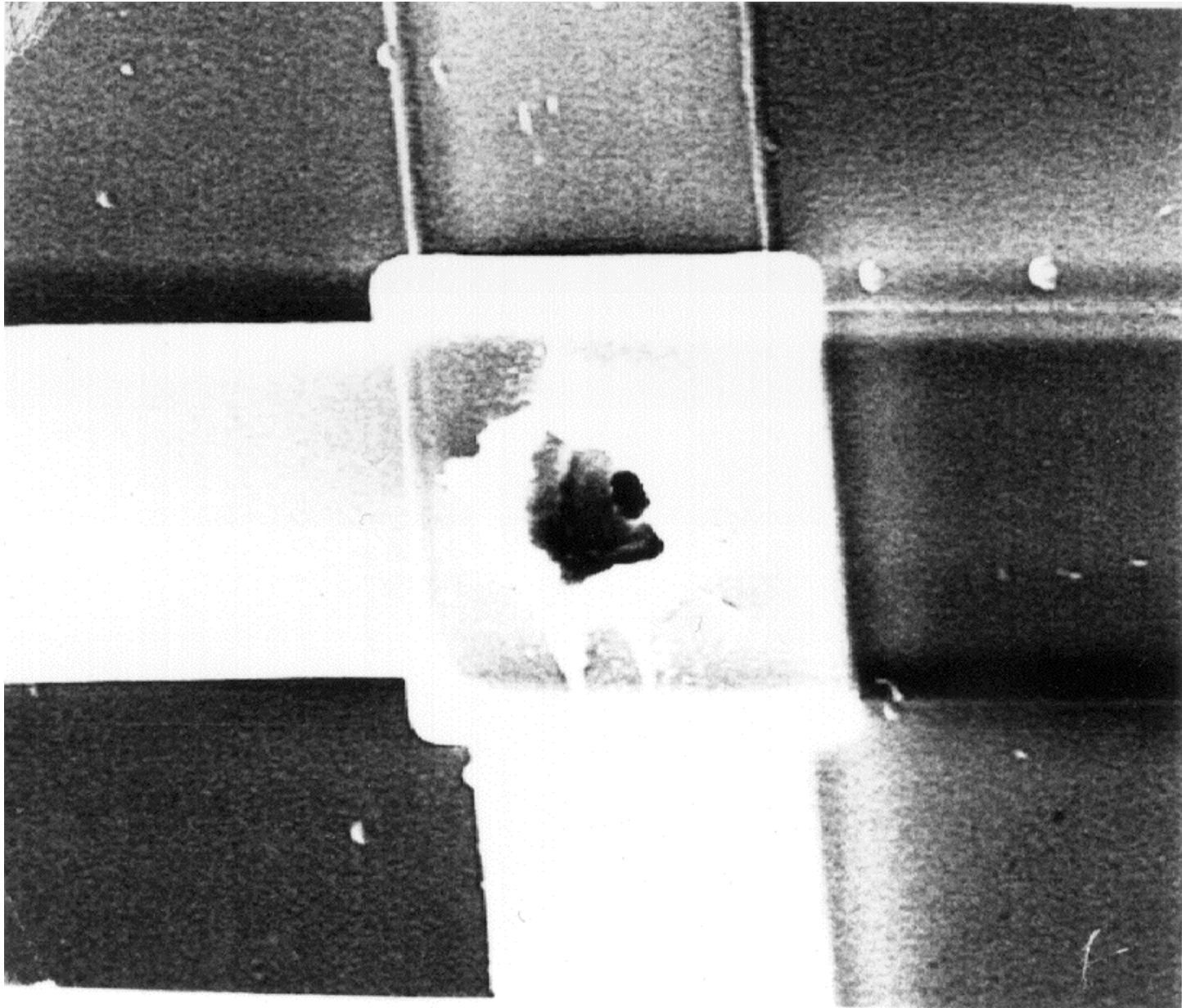
# Electromigration: Limits DC Current

Typical limit  
 **$1\text{mA}/\mu\text{m}$**



# Electromigration: **Need Multiple Contacts**

---



# Power and Clock Lines: Sizing

- Check for metal migration at worst power corner
- Do the following checks:

Process (np)	Temp	Voltage	Tests
<b>FF</b>	Low	High (min)	Power (DC), clock races, hold time
<b>SS</b>	High	Low (min)	Circuit speed, setup time
<b>SF</b>	Low	High (min)	Pseudo-nMOS noise margin, level shifters, memory write/read, ratioed circuits
<b>FS</b>	High	High (min)	Memories, ratioed circuits, level shifters

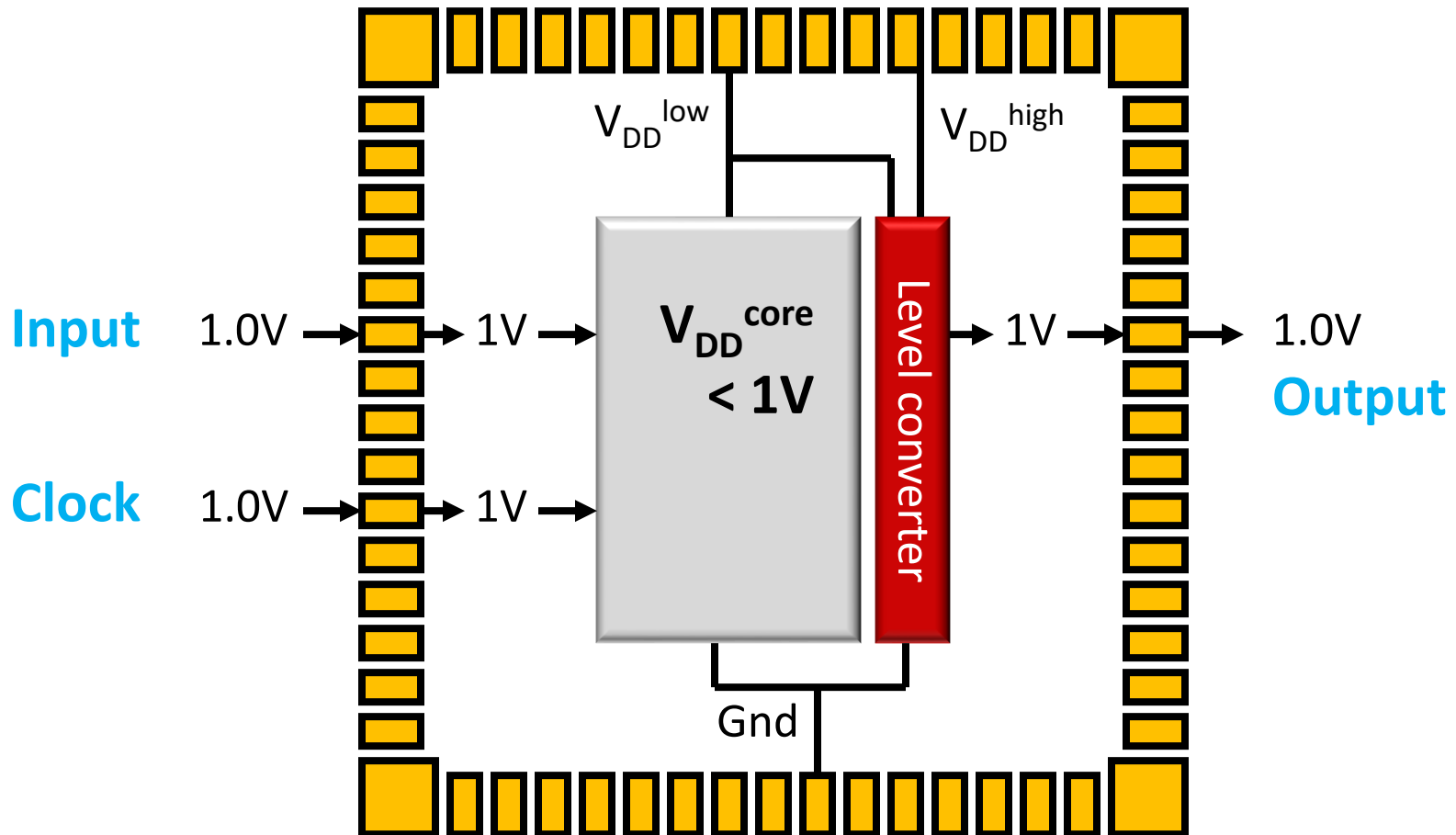
# Power Supply Rules

---

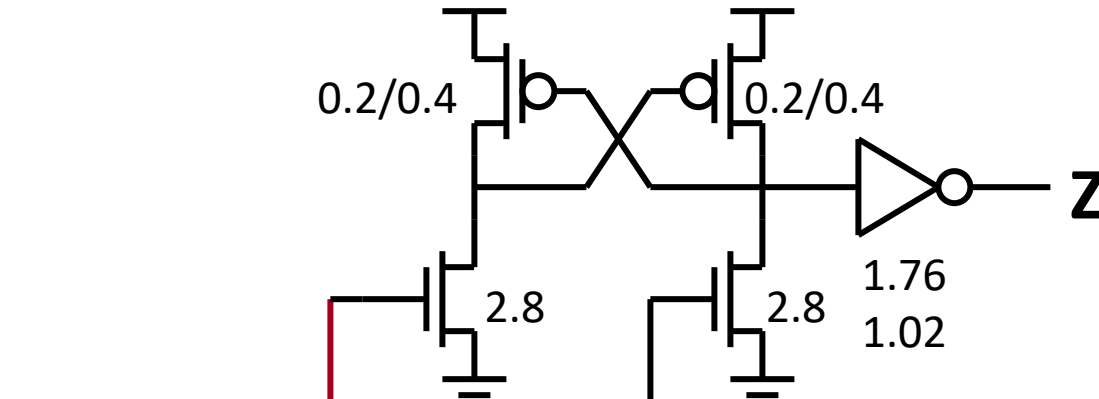
- **The exact rules depend on the technology**
  - Rules for resistance and electromigration
- **Example rules**
  - Ct for each  $16\lambda$  of transistor width (more is better)
  - Wire must have less than  $1\text{mA}/\mu\text{m}$  of width
  - Power/Gnd width =  $L_{\text{wire}} * \text{Sum (all transistors connected to wire)} / 3 * 10^6 \lambda$  (very approximate)
- **For small designs, supply design is less of an issue**
  - Total power is small
  - Chip is small, so wires are short
    - Will not be an issue in this class

# Working with Scaled $V_{DD}$

- Additional power/gnd pads, level shifter



# Simple On-chip Level Converter

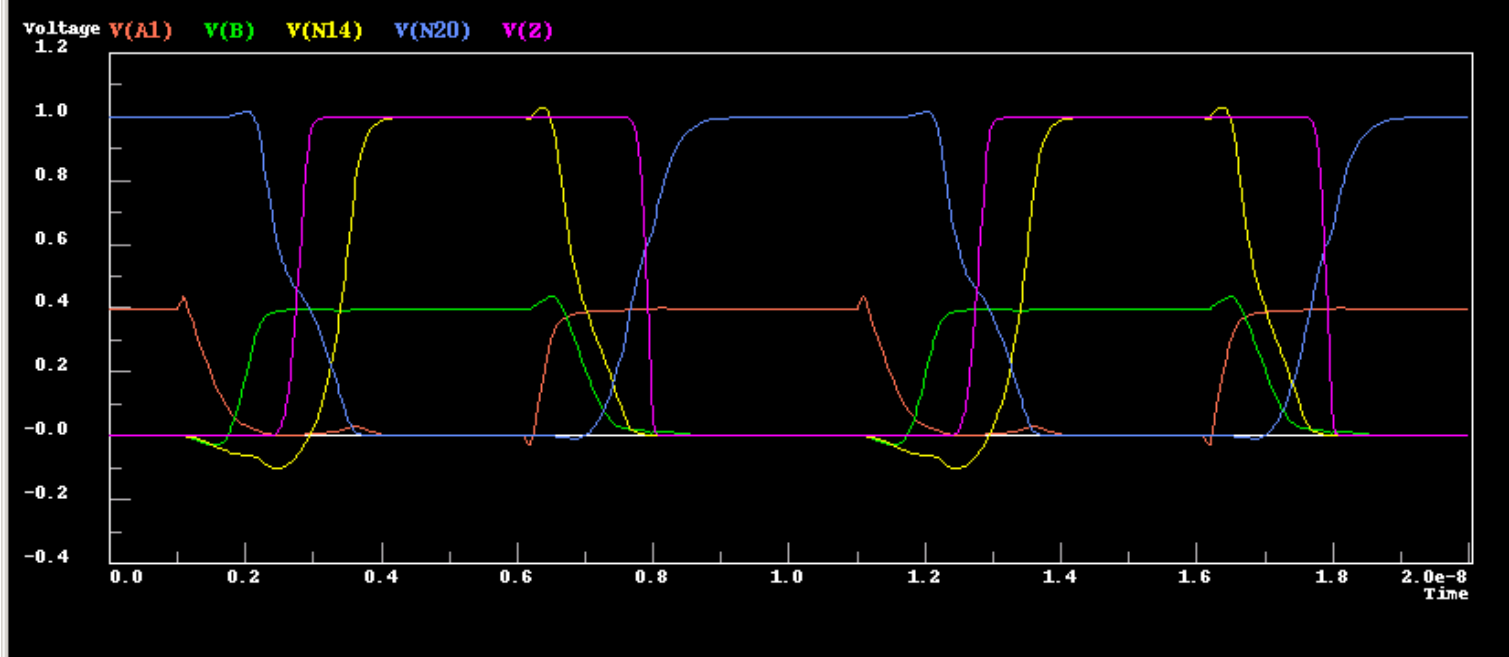
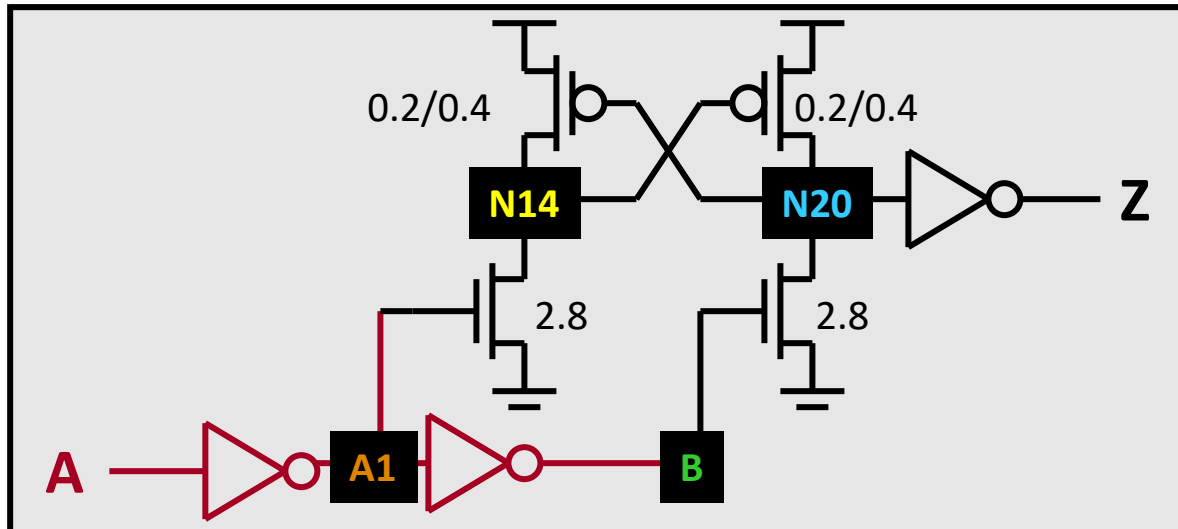


P: 1.76  
N: 1.02  
(4x inv)

1.76  
1.02

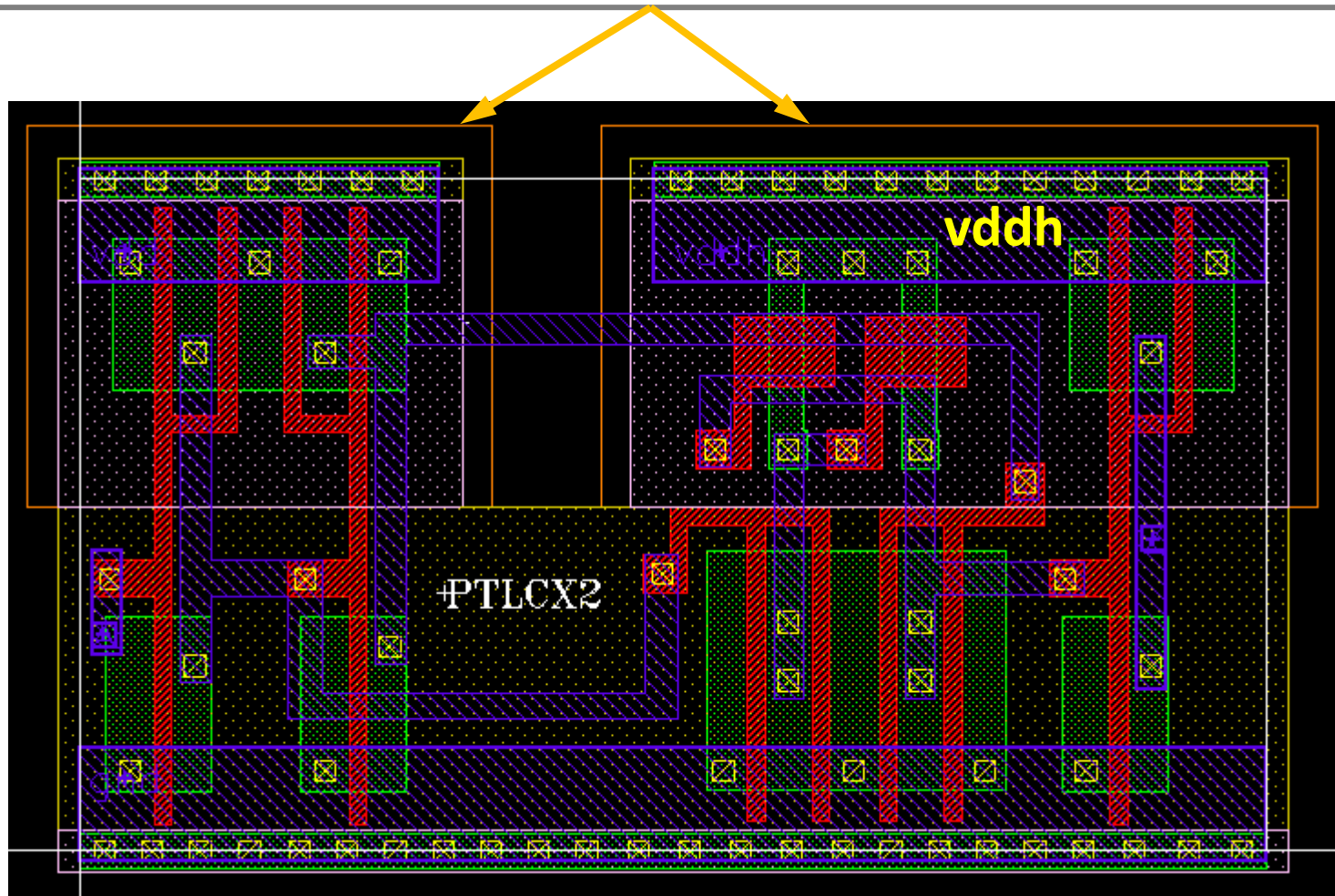
VddL (V)	tp:0-1 (ps)	tp:1-0 (ps)	tp:avg (ps)
0.4	1725	1750	1738
0.5	435	880	658
0.6	214	754	484
0.7	145	717	431
0.8	115	701	408
0.9	98	695	396
1.0	86	693	389

# Example: $V_{DDL} = 0.4V$ | 90nm





# Layout (Separate N-wells)



**PR boundary:  $7\mu\text{m} \times 3.9\mu\text{m}$  (std height)**

# Post-layout Simulation Results

---

VddL (V)	tp:0-1 (ps)	tp:1-0 (ps)	tp:avg (ps)	tp:avg (ps)
0.4	1725	1750	1738	1889
0.5	435	880	658	792
0.6	214	754	484	573
0.7	145	717	431	501
0.8	115	701	408	470
0.9	98	695	396	452
1.0	86	693	389	442

**Post-layout RC extraction** →

**Note:** final implementation may have additional buffers at the output to increase drive strength before the pads

- I/O pads: see W&H, 12.4.1-3